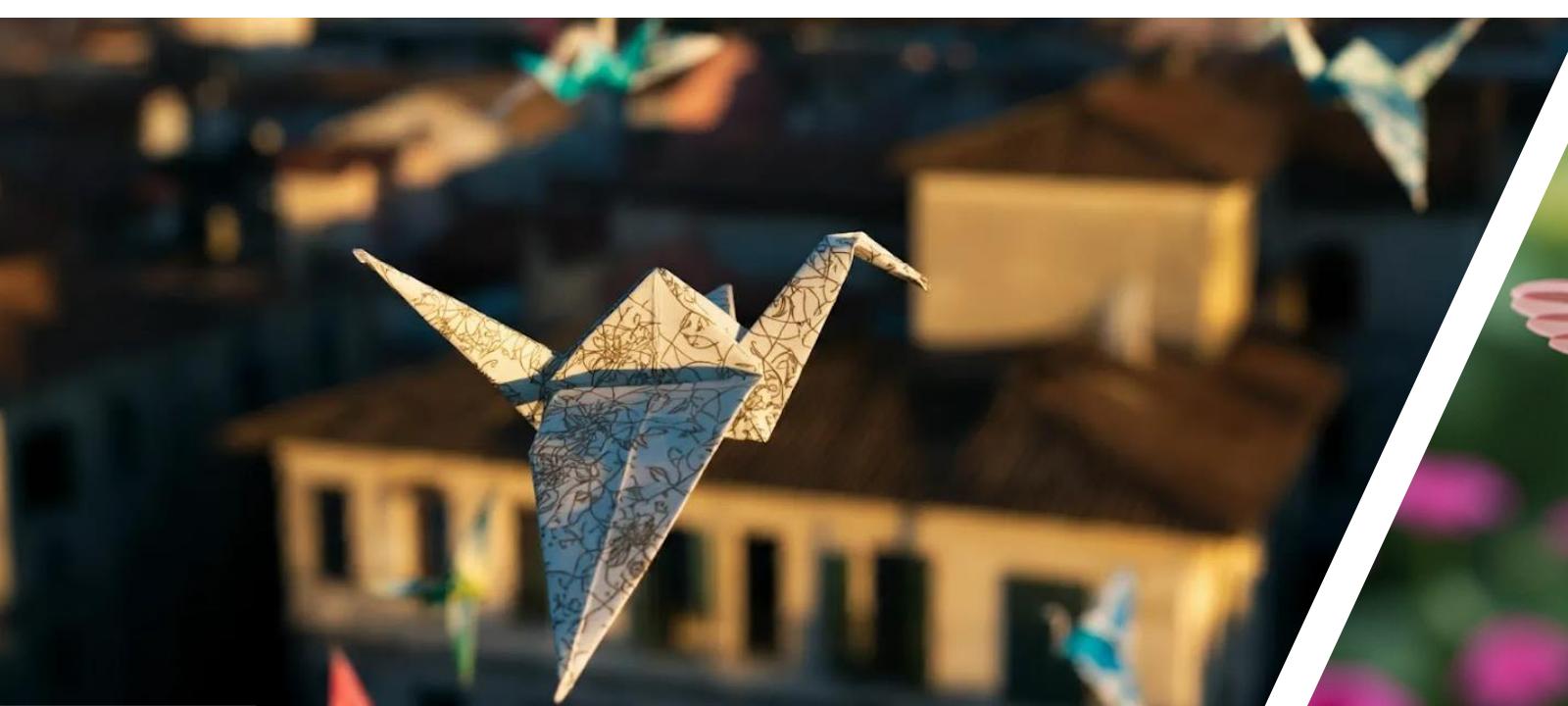
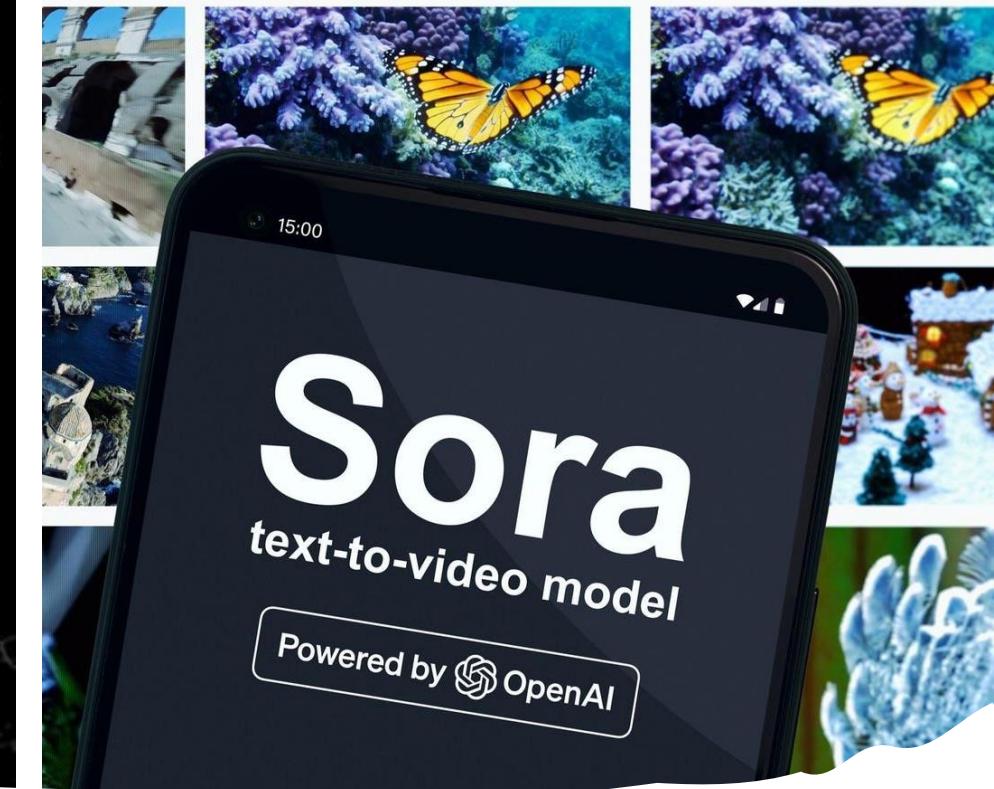


Diffusion Distillation in 2025

Tianwei Yin
Reve AI

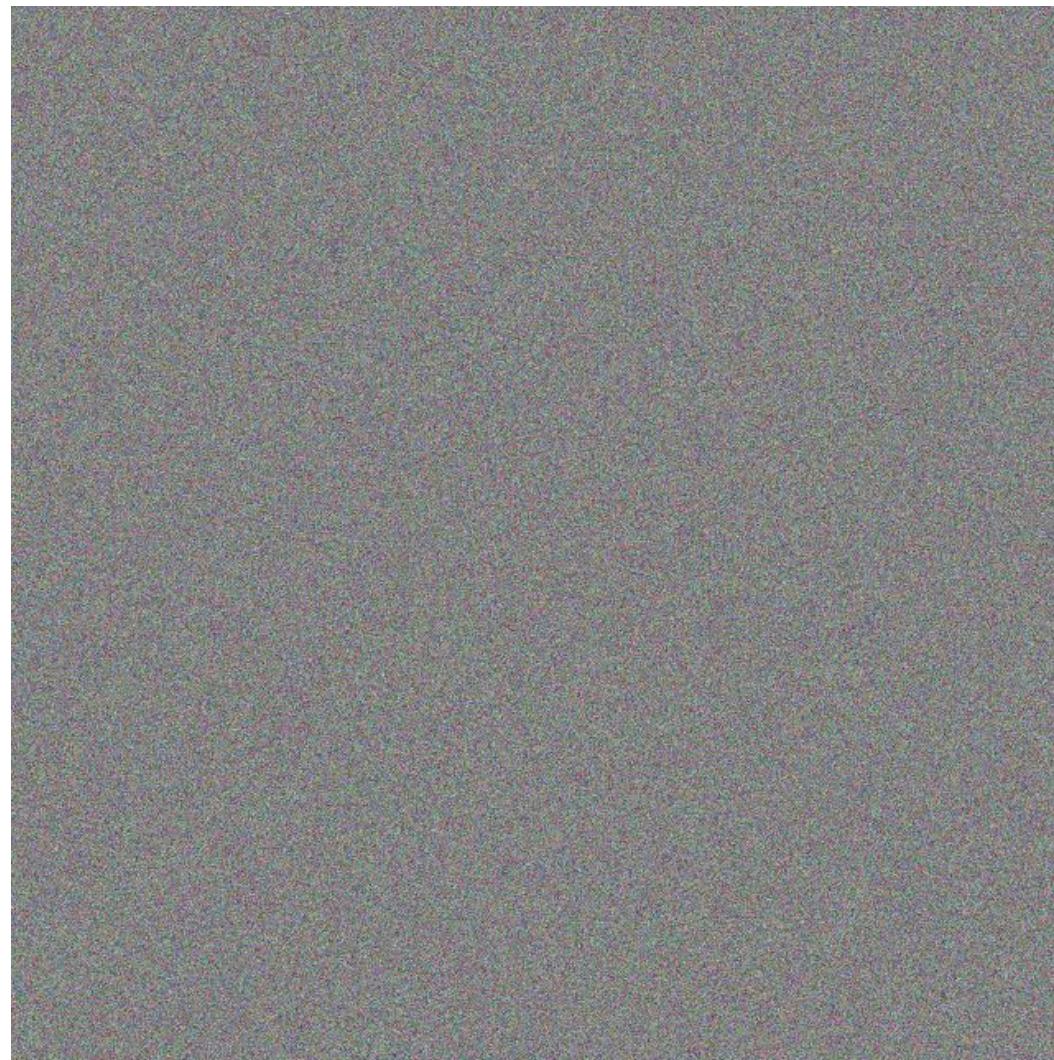


OpenAI launches Dall-E 3



Recent Advances in
Diffusion-based Generative
Models

Sampling with Diffusion Models is Slow



- Orders of magnitude slower than GANs and VAEs
- Iterative denoising process
 - 20-50 steps for reasonable images
- We speed it up by step distillation!



Overview of Diffusion Distillation

Trajectory Preserving

Knowledge Distillation [LL, 2020]

Progressive Distillation [SH, ICLR 2022]

Guidance Distillation [MRGKEHS, CVPR 2023]

Consistency Models [SDCS, ICML 2023]

Distribution Matching

DMD [YGZSDFP, CVPR 2024]

DMD 2 [YGPZSDF, NeurIPS 2024]

Diff-Instruct [LHZSLZ, NeurIPS 2023]

MMD [SMHH, NeurIPS 2024]

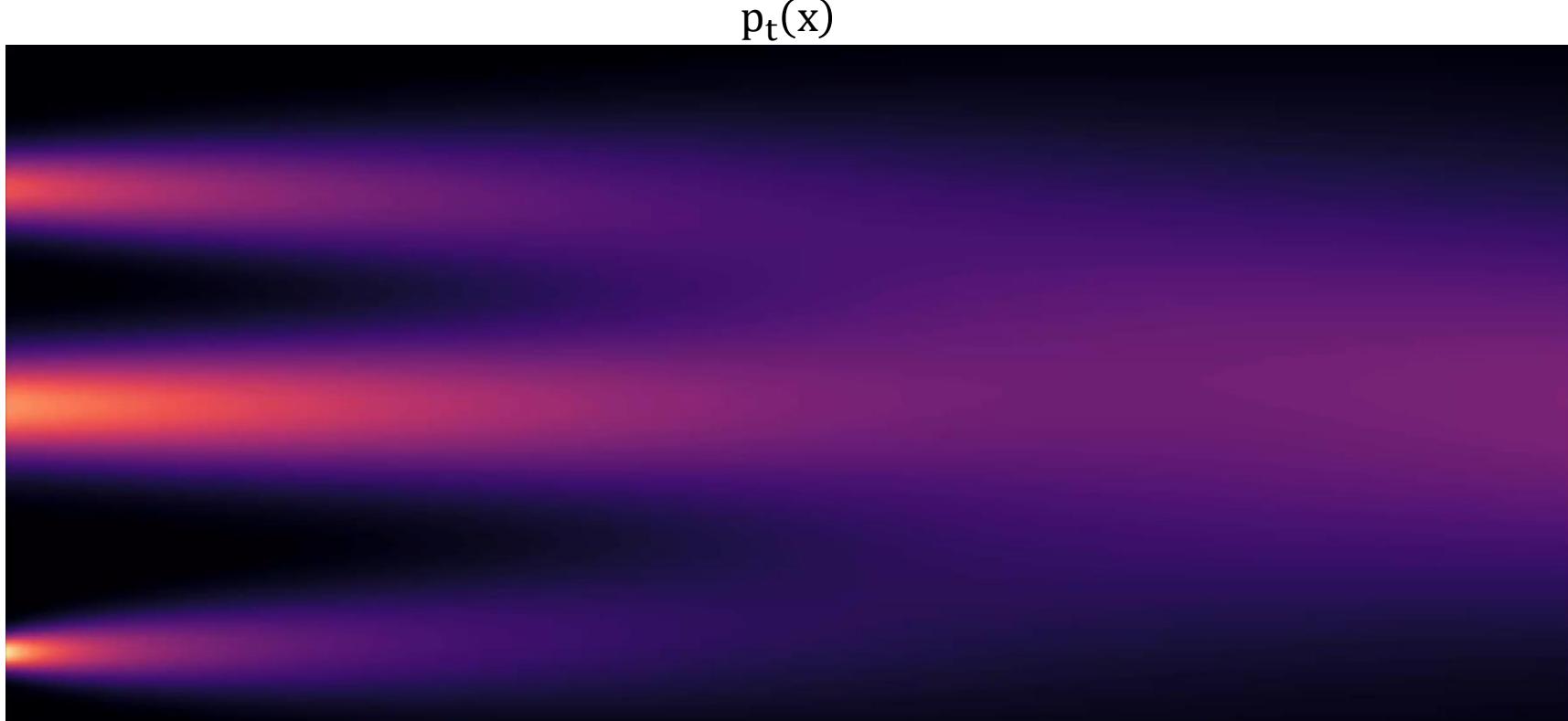
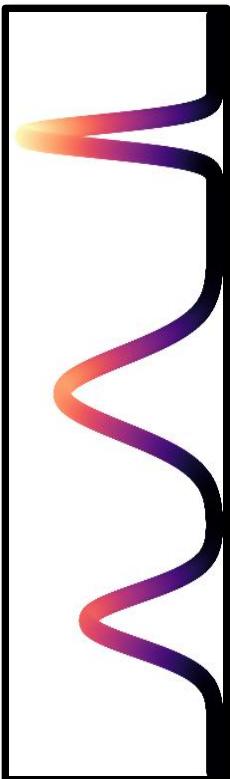
CausVid [YZZFDSH, CVPR 2025]

GAN

Too many, covered in Junyan's talk.

Diffusion Basics

Converting Data Distribution to Gaussian

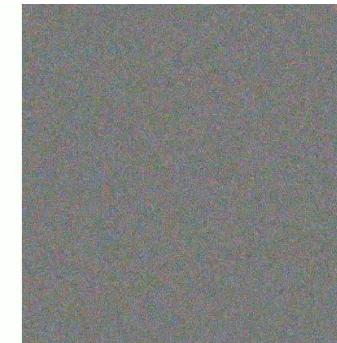
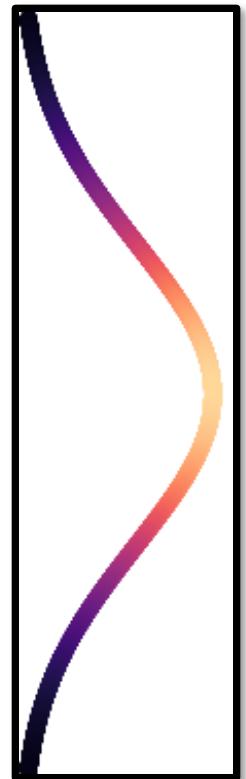
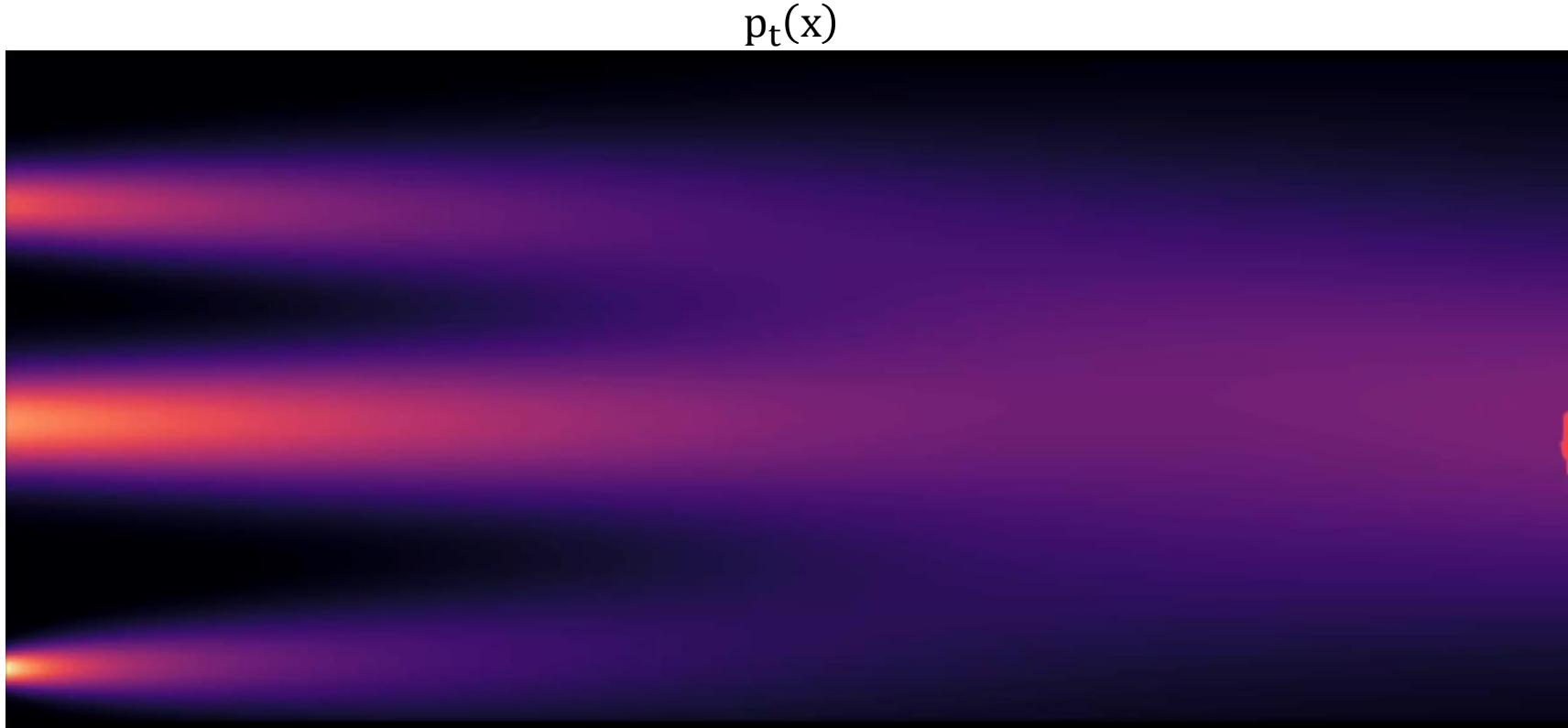
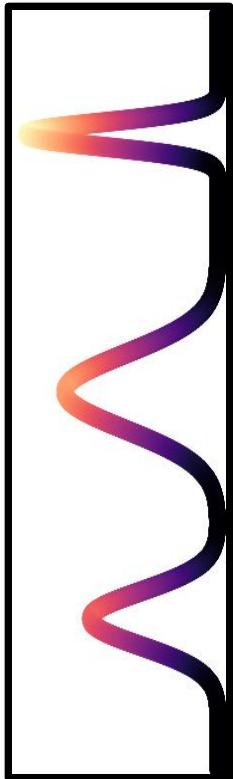


$$\begin{aligned} p_0(x) \\ = \\ p_{\text{data}}(x) \end{aligned}$$



$$\begin{aligned} p_T(x) \\ \approx \\ \mathcal{N}(0, 1) \end{aligned}$$

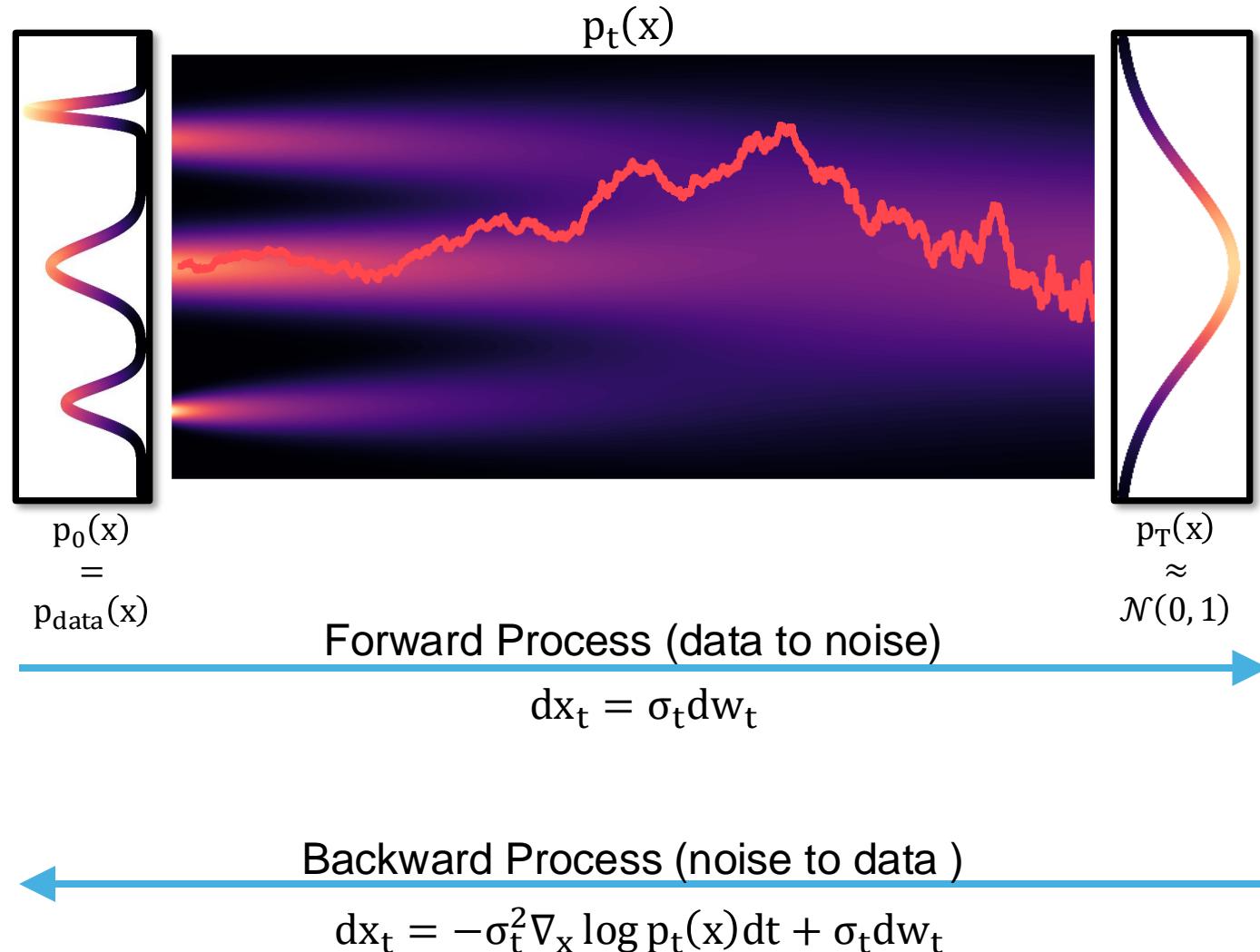
Generating Sample from Noise



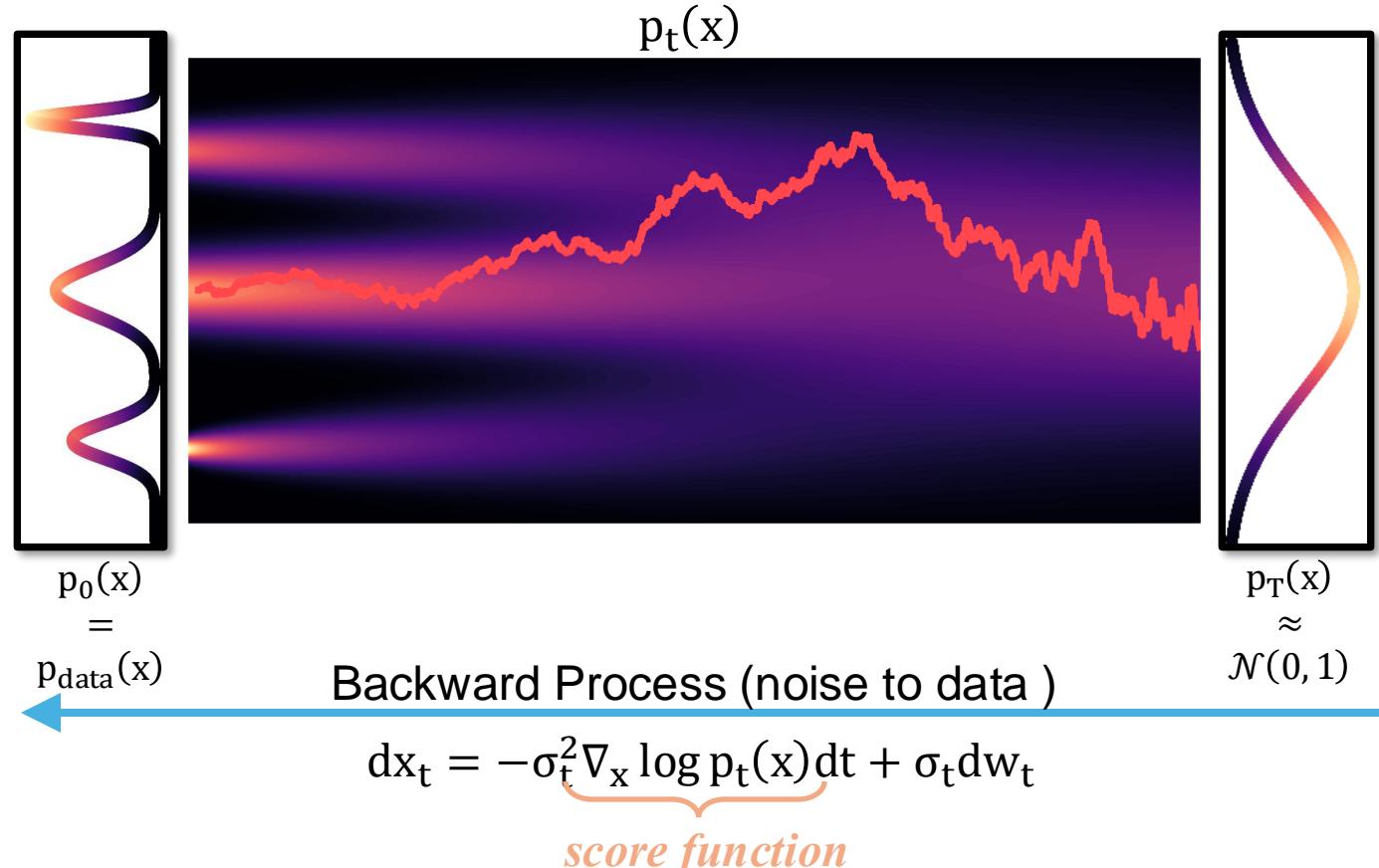
$$\begin{aligned} p_0(x) \\ = \\ p_{\text{data}}(x) \end{aligned}$$

$$\begin{aligned} p_T(x) \\ \approx \\ \mathcal{N}(0, 1) \end{aligned}$$

Basics of Diffusion Model



Basics of Diffusion Model

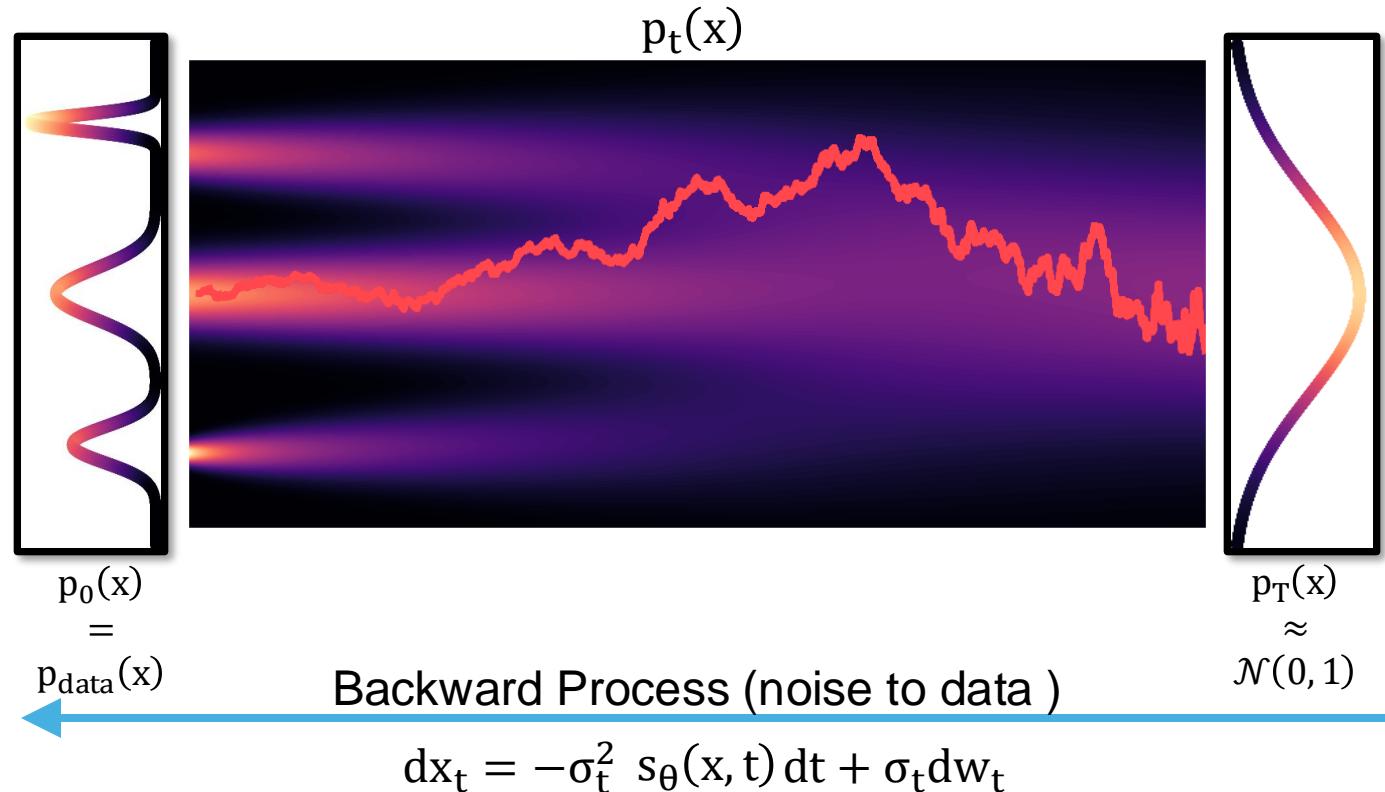


$$\underbrace{s_\theta(x, t)}_{\text{Neural Network}} = \nabla_x \log p_t(x)$$

Score Matching Loss:

$$\mathbb{E}_{t,x}[\lambda(t) \|\nabla_x \log p_t(x) - s_\theta(x, t)\|]$$

Basics of Diffusion Model



If we omit the intermediate noise injection ($\sigma_t dw_t$), we obtain a backward ODE that establishes one-to-one mappings between noise and data.

We can utilize any ODE solvers (DDIM, Euler, etc.)

Overview of Diffusion Distillation

Trajectory
Preserving

Knowledge Distillation [LL, 2020]

Progressive Distillation [SH, ICLR 2022]

Guidance Distillation [MRGKEHS, CVPR 2023]

Consistency Models [SDCS, ICML 2023]

Distribution
Matching

DMD [YGZSDFP, CVPR 2024]

DMD 2 [YGPZSDF, NeurIPS 2024]

Diff-Instruct [LHZSLZ, NeurIPS 2023]

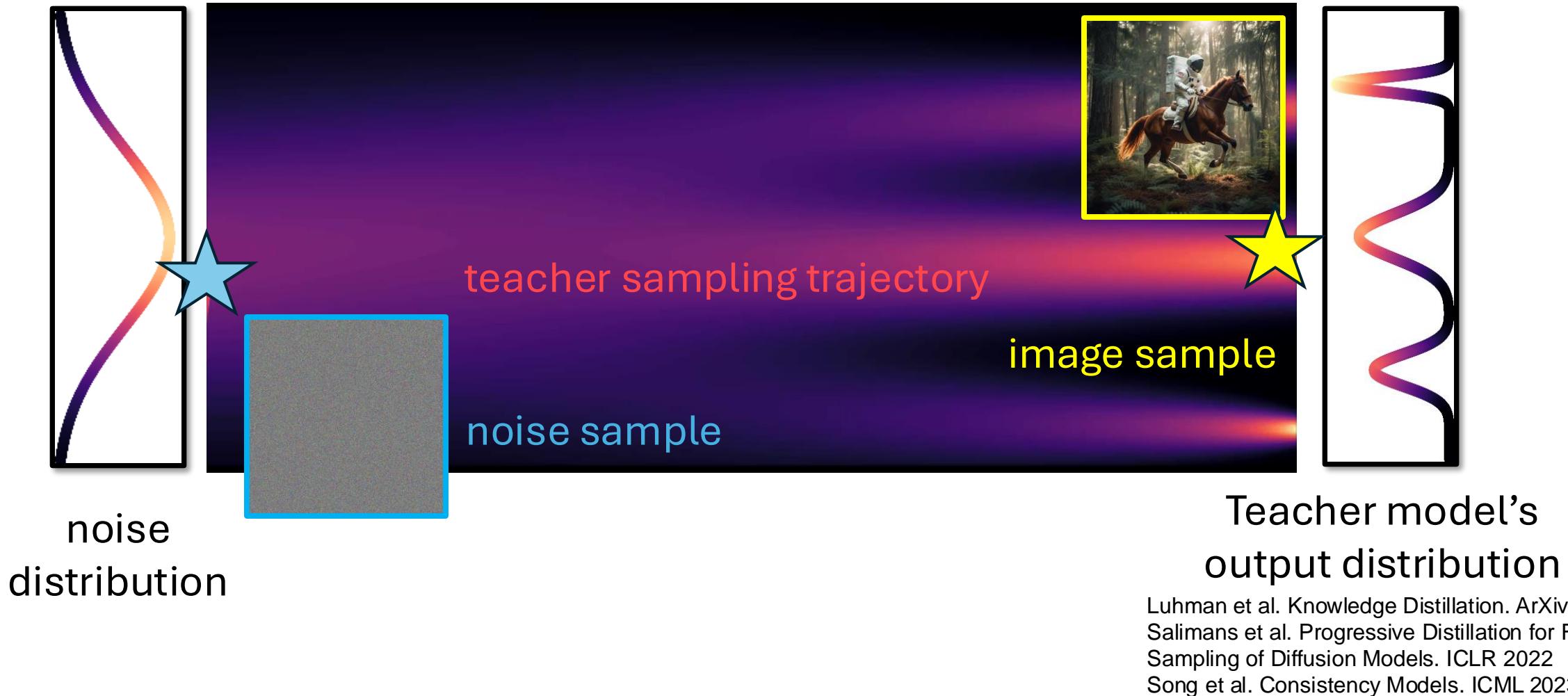
MMD [SMHH, NeurIPS 2024]

CausVid [YZZFDSP, CVPR 2025]

GAN

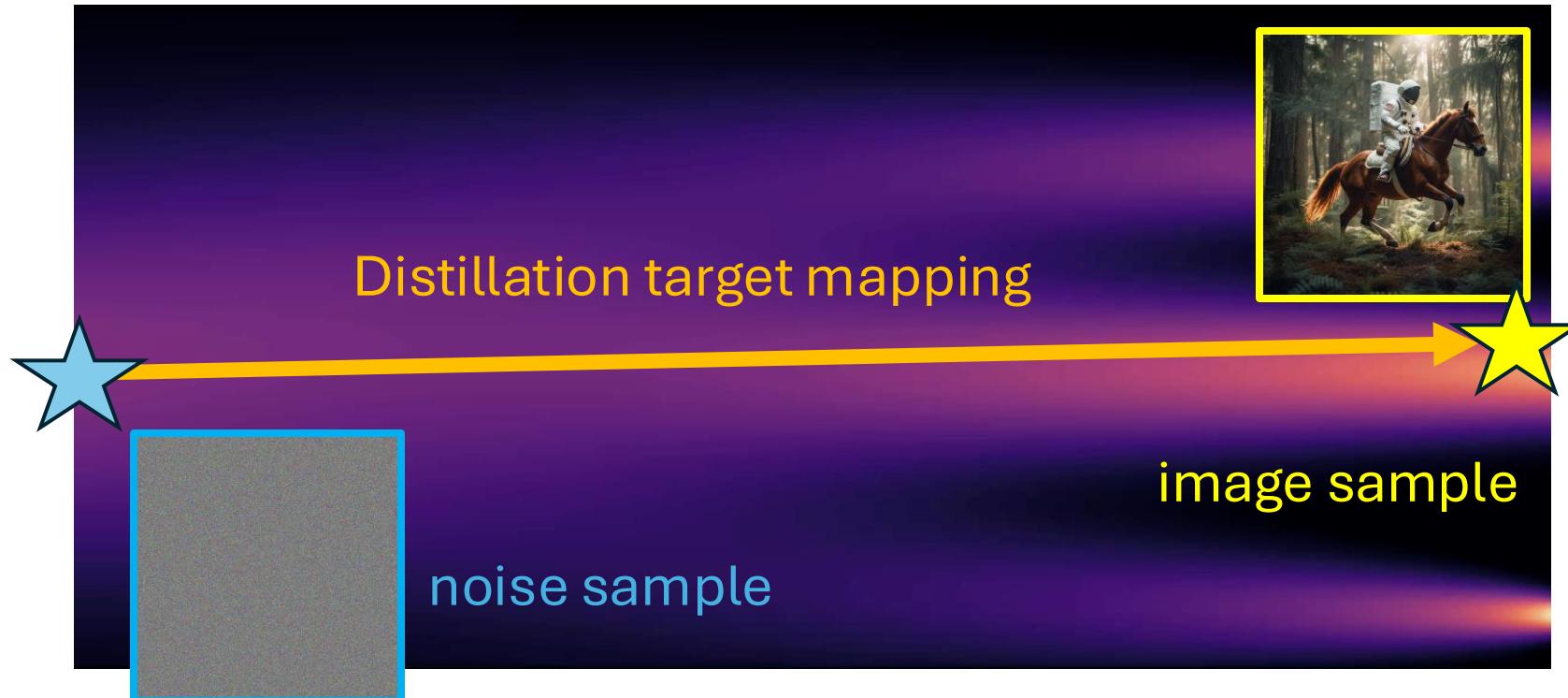
Too many, covered in Junyan's talk.

Diffusion distillation, its challenges

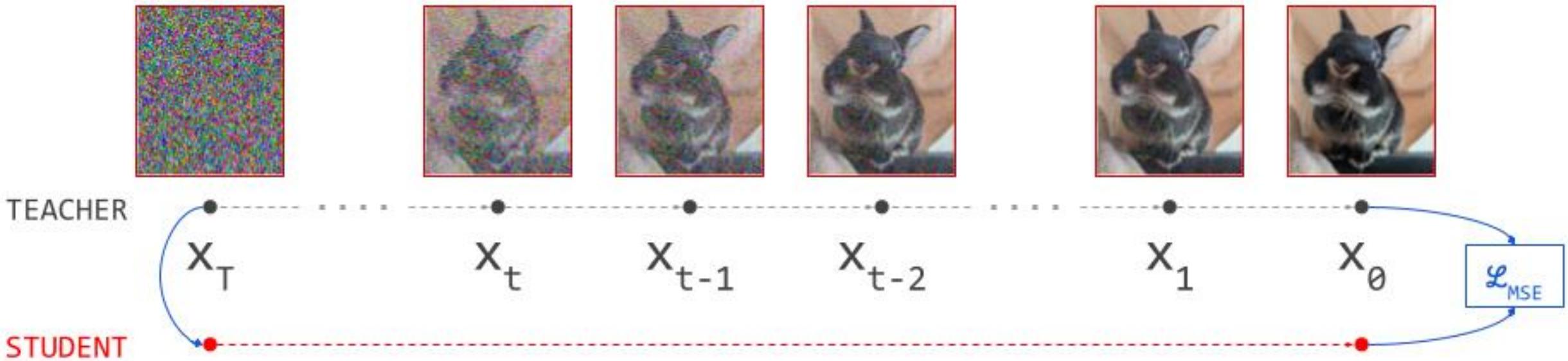


Luhman et al. Knowledge Distillation. ArXiv 2021
Salimans et al. Progressive Distillation for Fast Sampling of Diffusion Models. ICLR 2022
Song et al. Consistency Models. ICML 2023

Trajectory-preserving distillation tries to recover **pointwise mapping** from noise-to-image



Knowledge Distillation: Distilling diffusion sampling into a single forward pass



Distill the entire sampling procedure into a network with the same architecture used for a single diffusion prediction step, by matching outputs with a MSE loss

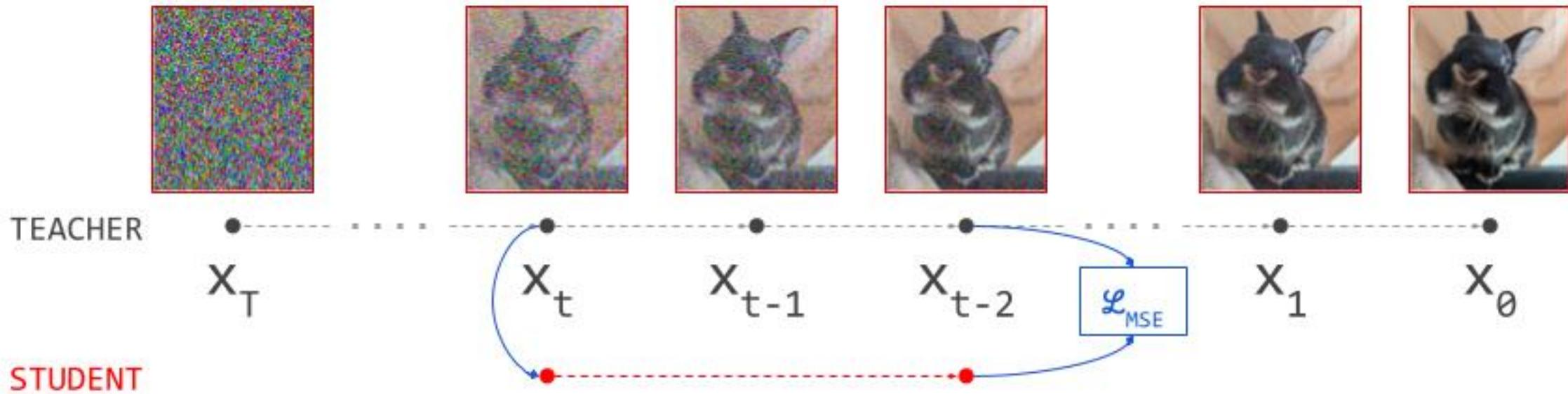
Proof-of-Concept: We really don't need 50 or 1000 steps!

Slow to Run: Requires full trajectory sampling

Low Performance Upper Bound: Error accumulation in deterministic sampling

Luhman et al. Knowledge Distillation.
ArXiv 2021
Image credit:
<https://sander.ai/2024/02/28/paradox.html>

Progressive Distillation: Iteratively halve the number of sampling steps



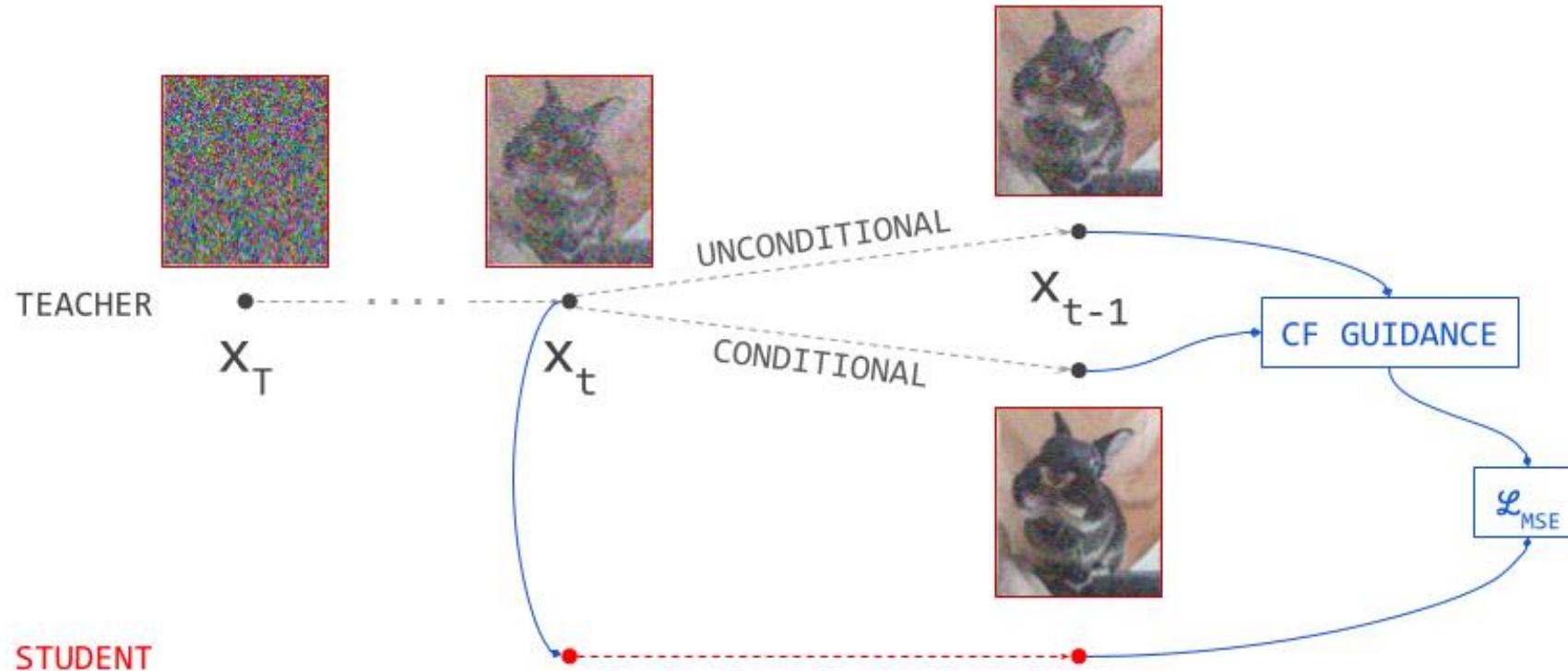
The teacher uses the full diffusion process (e.g., many time steps). The student then learns to produce similar outputs but in half the number of steps.

Faster-to-Run: No precomputation needed!

Better quality tradeoff: 8 step generally leads to satisfactory result

Salimans, Tim, and Jonathan Ho. Progressive distillation for fast sampling of diffusion models.
Image credit:
<https://sander.ai/2024/02/28/paradox.html>

Guidance Distillation: Avoid two network passes in classifier-free guidance



Approximate guided score with a single network forward pass.

Effective and no quality drop.
Used in Flux-dev image generation

Meng, Chenlin, et al. "On distillation of guided diffusion models." CVPR 2023.
Image credit:
<https://sander.ai/2024/02/28/paradox.html>

Overview of Diffusion Distillation

Trajectory Preserving

Knowledge Distillation [LL, 2020]

Progressive Distillation [SH, ICLR 2022]

Guidance Distillation [MRGKEHS, CVPR 2023]

Consistency Models [SDCS, ICML 2023]

Distribution Matching

DMD [**YGZSDFP**, CVPR 2024]

DMD 2 [**YGPZSDF**, NeurIPS 2024]

Diff-Instruct [LHZSLZ, NeurIPS 2023]

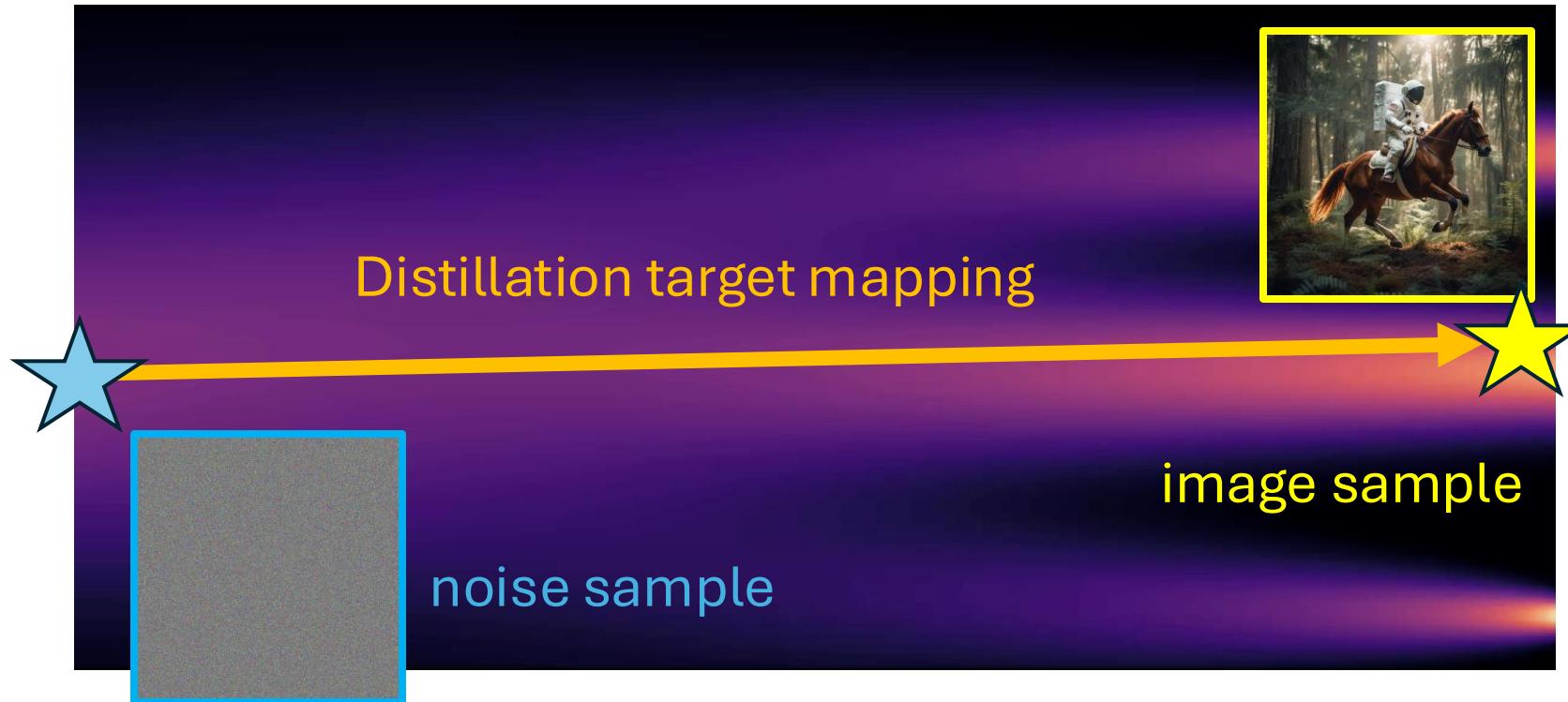
MMD [SMHH, NeurIPS 2024]

CausVid [**YZZFDSH**, CVPR 2025]

GAN

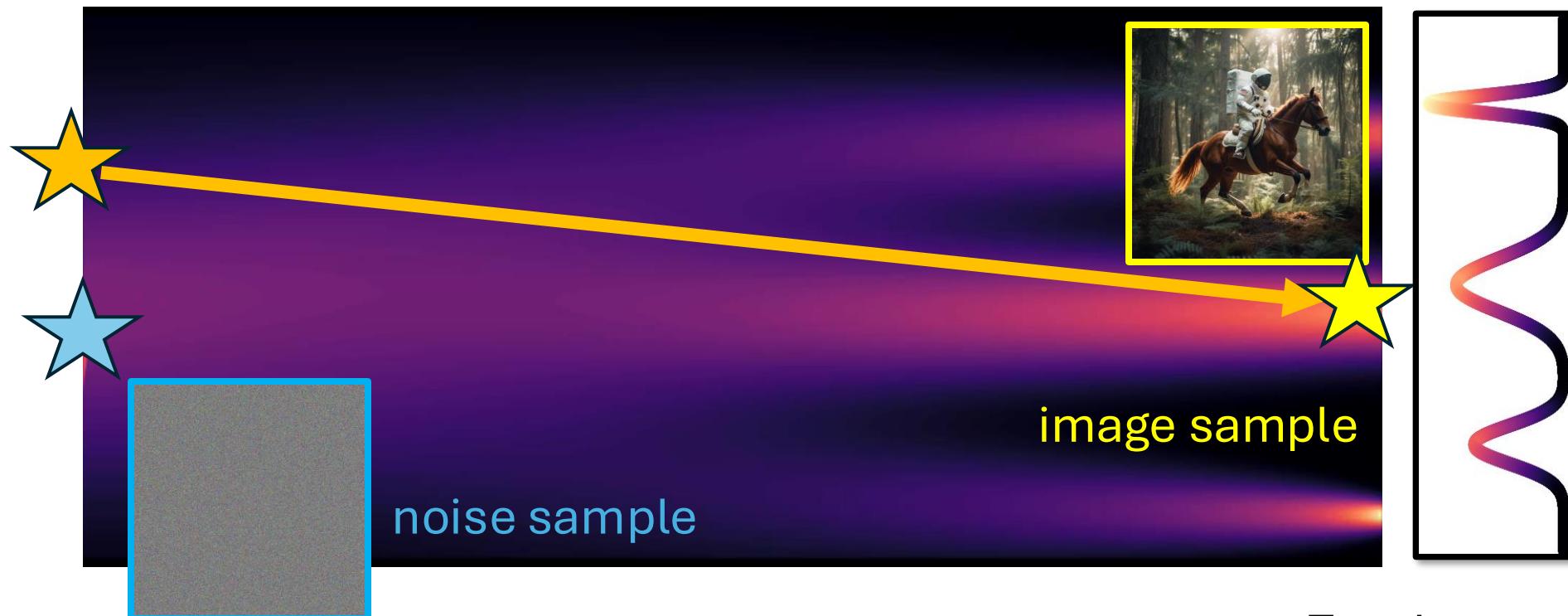
Too many, covered in Junyan's talk.

Trajectory-preserving distillation tries to recover **pointwise mapping** from noise-to-image



DMD matches the teacher's distribution

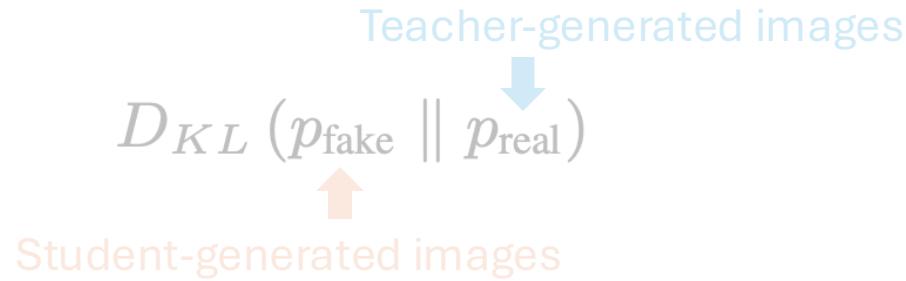
learning the exact noise -> image mapping is unnecessarily hard



Different Mapping, Same Quality, Easier to train

Teacher model's
output distribution

KL Divergence Minimization using Diffusion Models



$$\nabla_{\theta} D_{KL} = \mathbb{E}_{\substack{z \sim \mathcal{N}(0; \mathbf{I}) \\ x = G_{\theta}(z)}} \left[- (s_{\text{real}}(x) - s_{\text{fake}}(x)) \nabla_{\theta} G_{\theta}(z) \right]$$

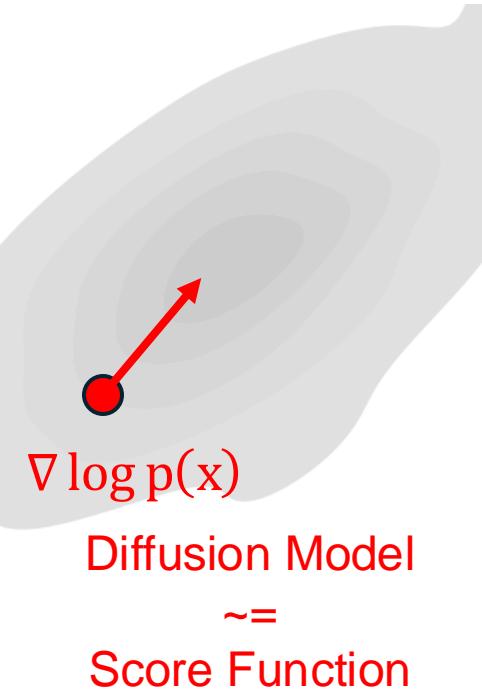
where $s_{\text{real}}(x) = \nabla_x \log p_{\text{real}}(x)$, $s_{\text{fake}}(x) = \nabla_x \log p_{\text{fake}}(x)$

Key insight

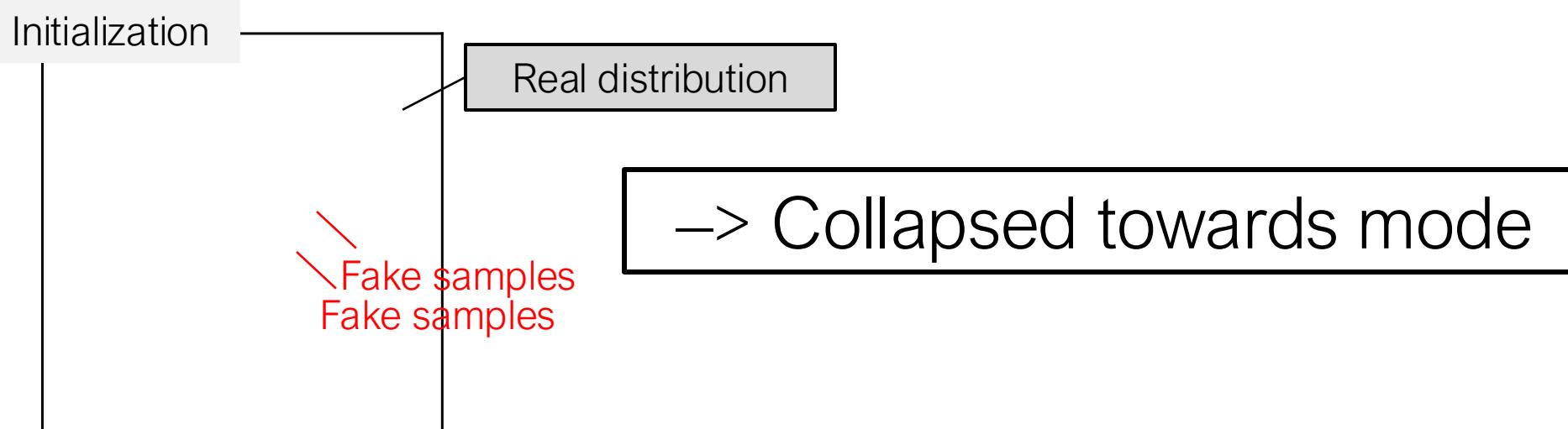
We can approximate these two gradients using 2 diffusion models!

$$\nabla_{\theta} D_{KL} = \mathbb{E}_{\substack{z \sim \mathcal{N}(0; \mathbf{I}) \\ x = G_{\theta}(z)}} \left[- (s_{\text{real}}(x) - s_{\text{fake}}(x)) \nabla_{\theta} G_{\theta}(z) \right]$$

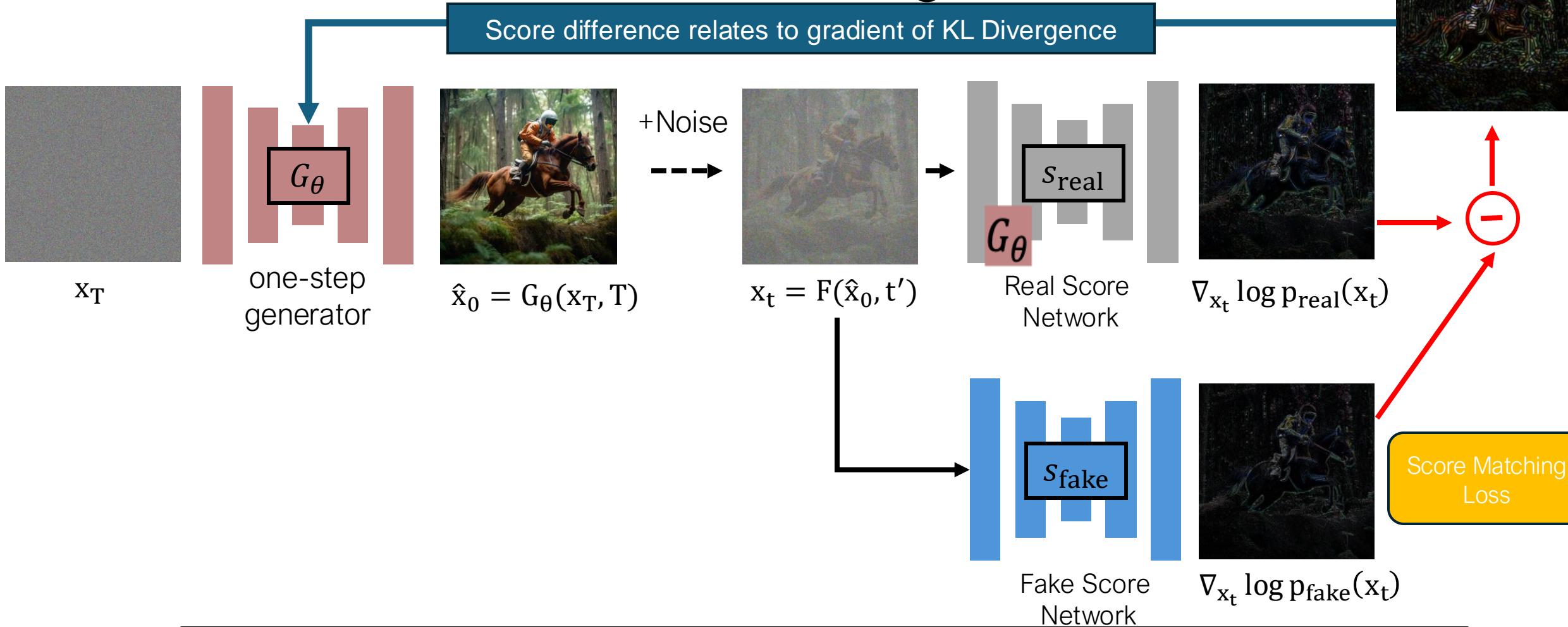
where $s_{\text{real}}(x) = \nabla_x \log p_{\text{real}}(x)$, $s_{\text{fake}}(x) = \nabla_x \log p_{\text{fake}}(x)$



Can we supervise a 1-step Generator using a single pretrained Score Function?

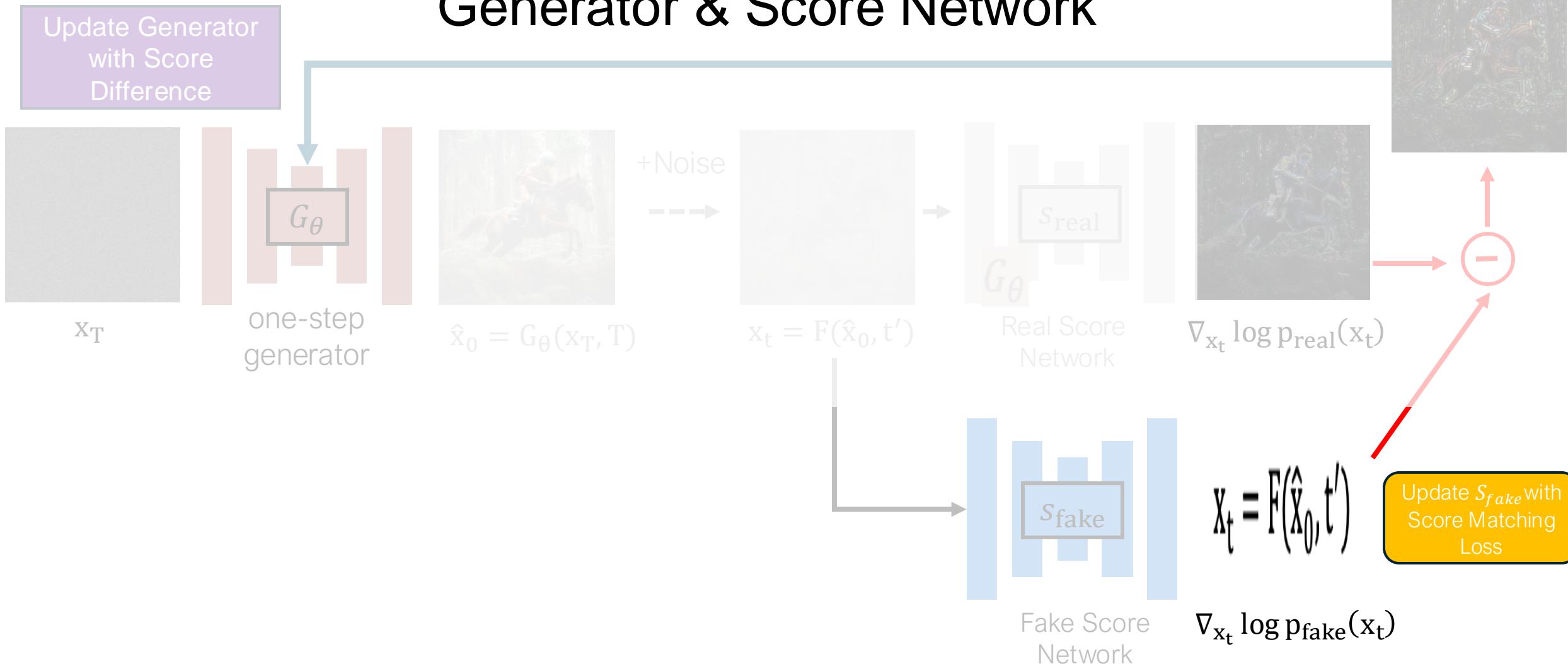


Training a second Score Function to estimate the KL gradient

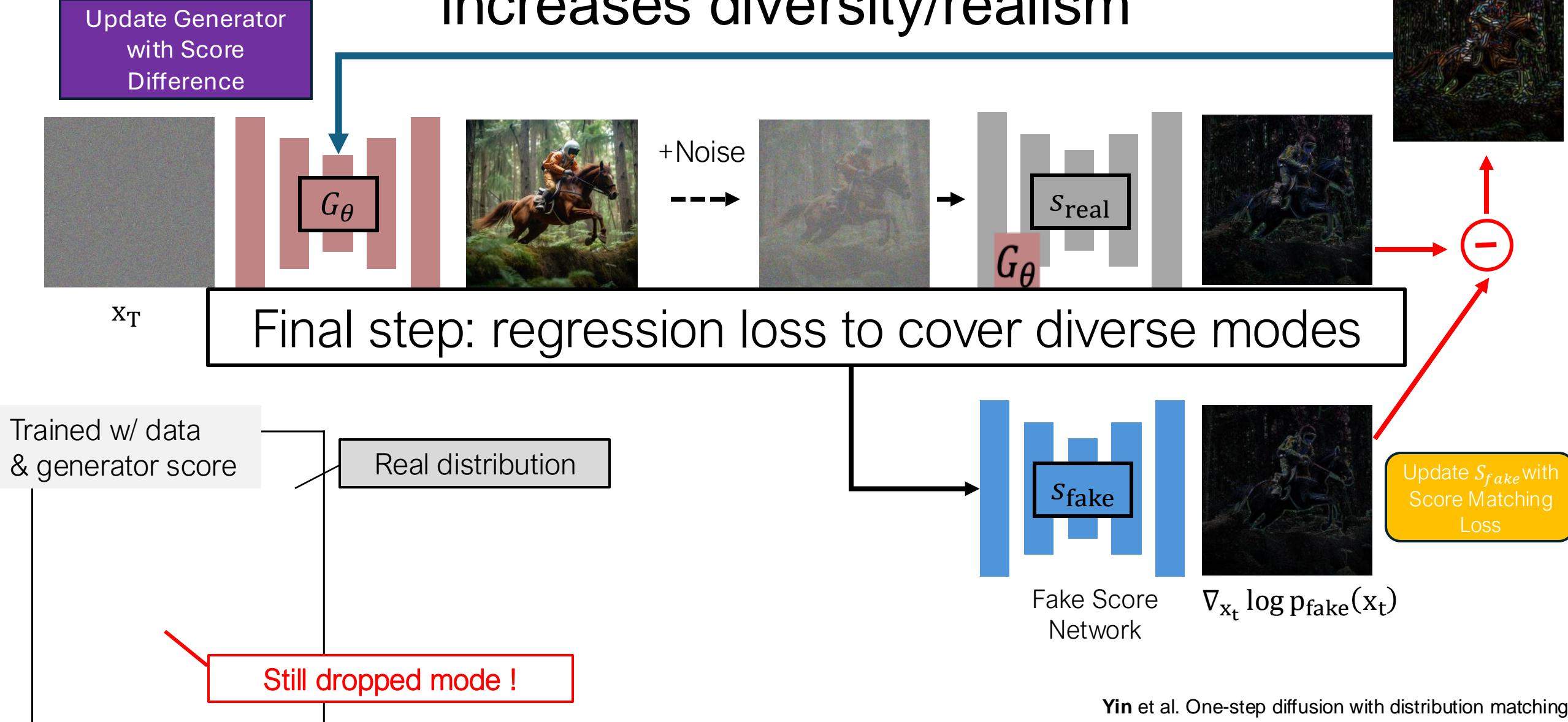


Train a score function on the *fake/generated* images

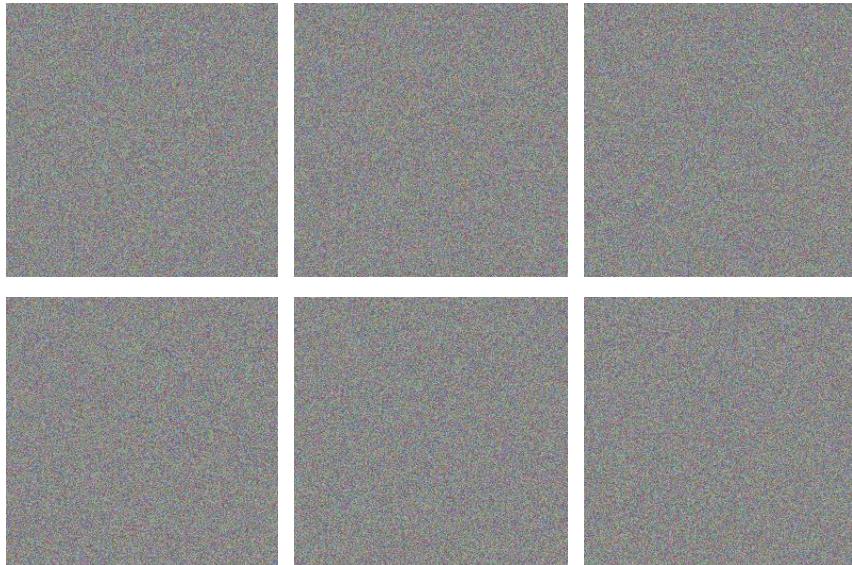
Alternatively optimize Generator & Score Network



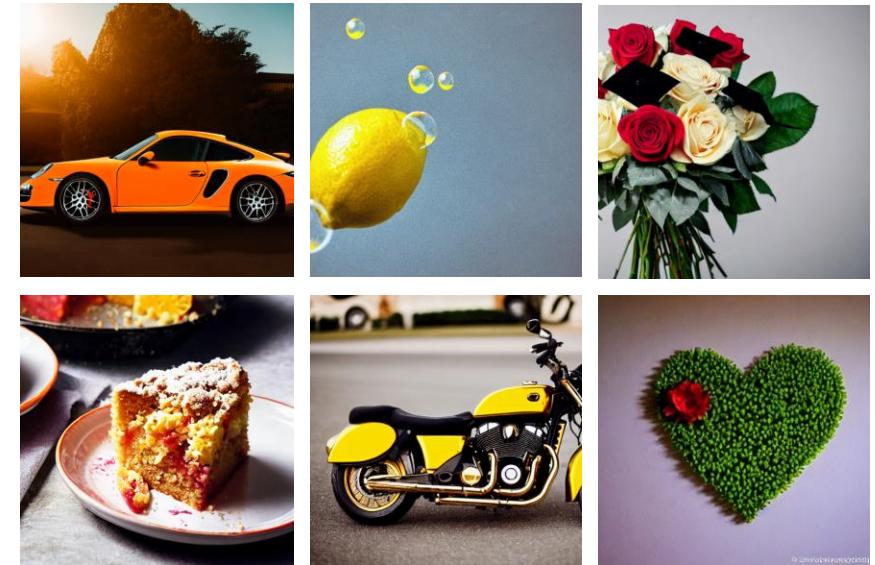
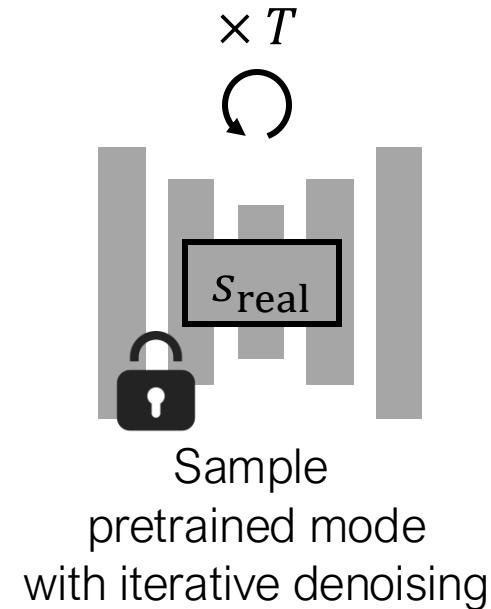
Distribution matching increases diversity/realism



Noise-Image Pair Collection

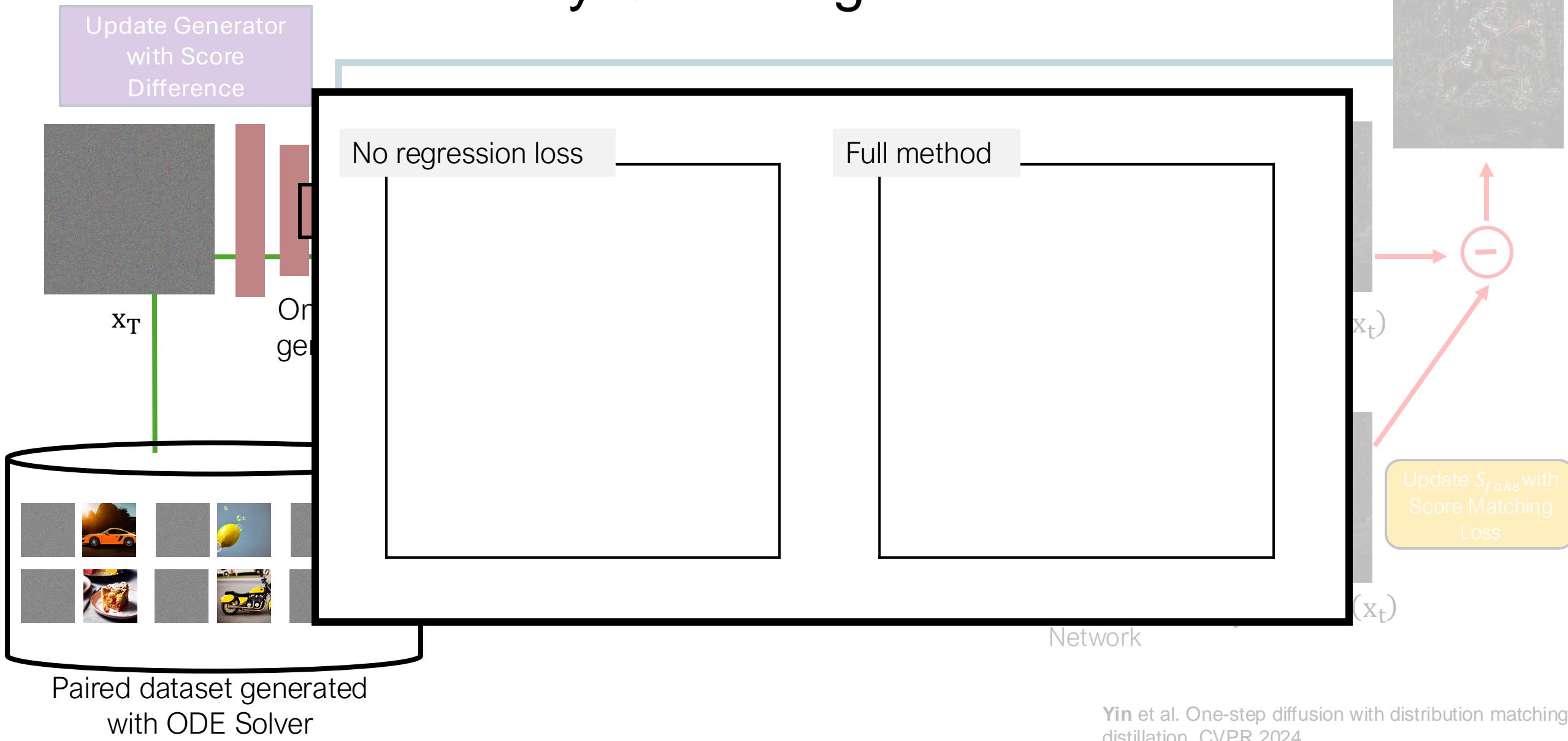


input noise maps



output images

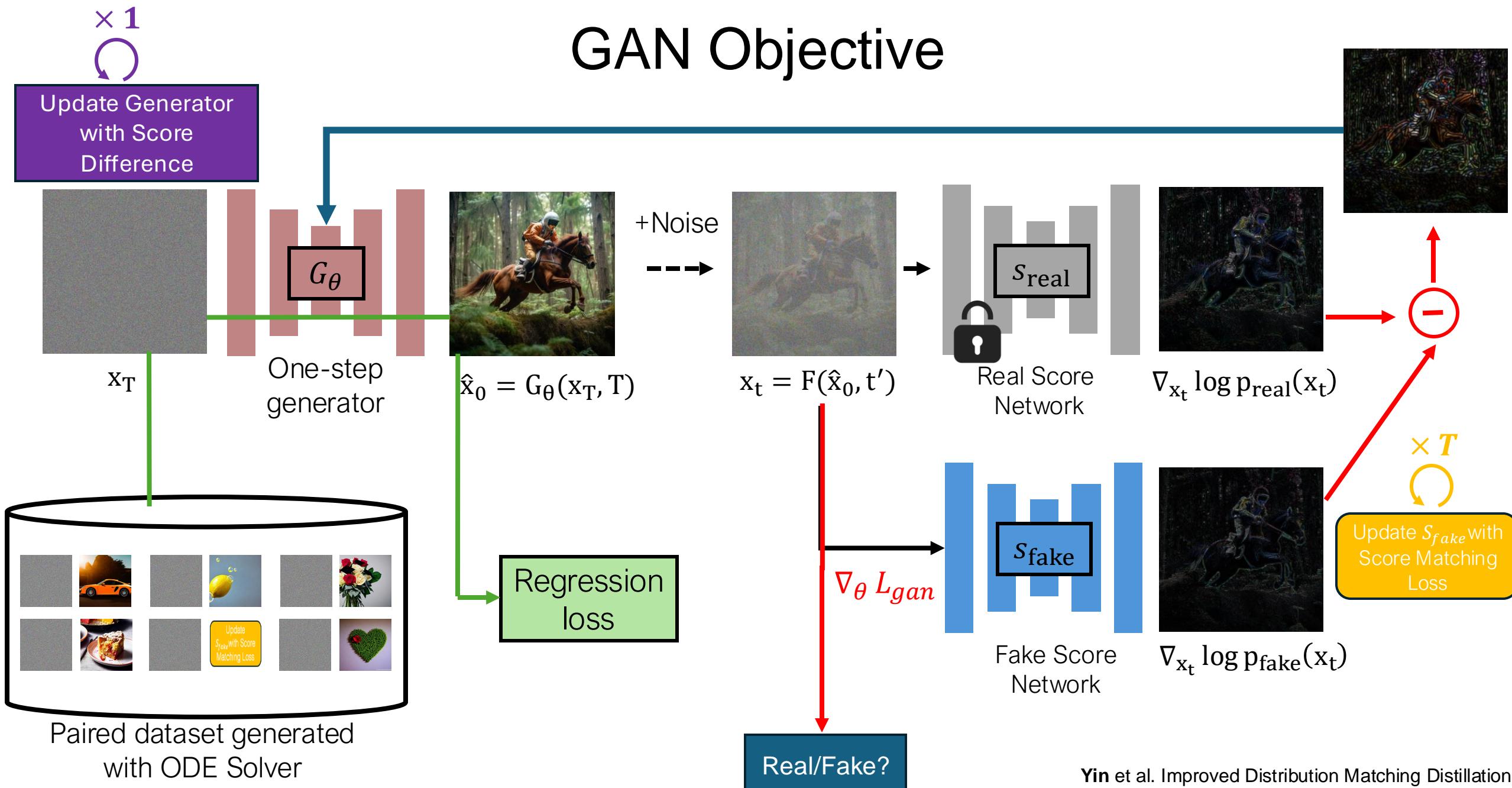
Auxiliary ODE Regression Loss



Improvement 1

GAN

GAN Objective





No GAN Loss



Ours w. GAN Loss

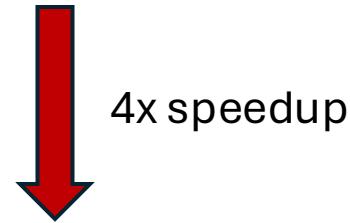
Improvement 2

Multi-step Generation

With no domain shift & quality loss, minimal additional latency

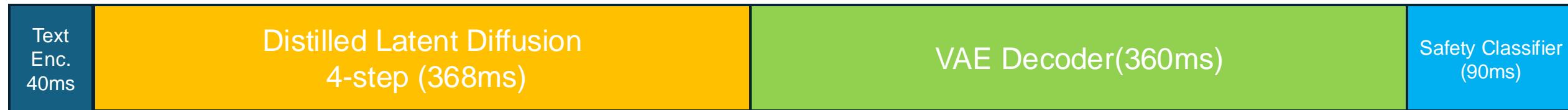
Few-steps diffusion latency

Distilled Latent Diffusion
4-step (368ms)



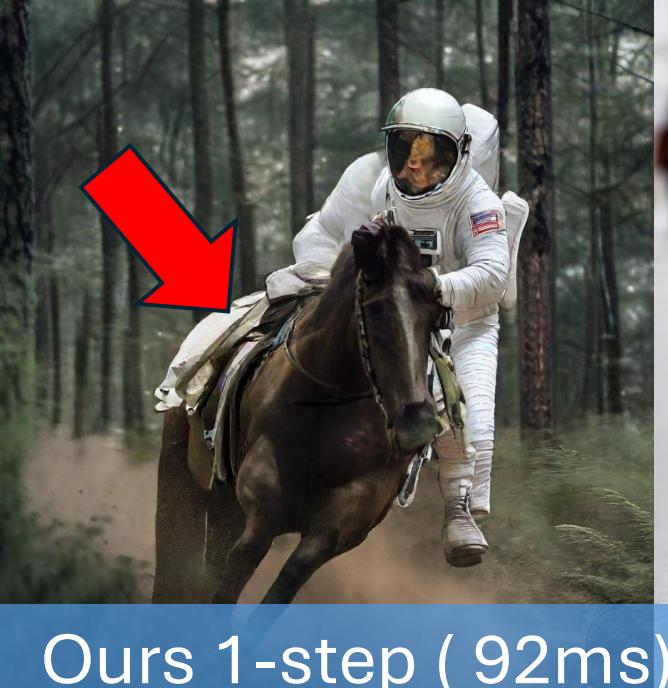
Distilled
Latent Diffusion
1-step (92ms)

Text2Image Full Pipeline Latency Breakdown



Only 30% faster!





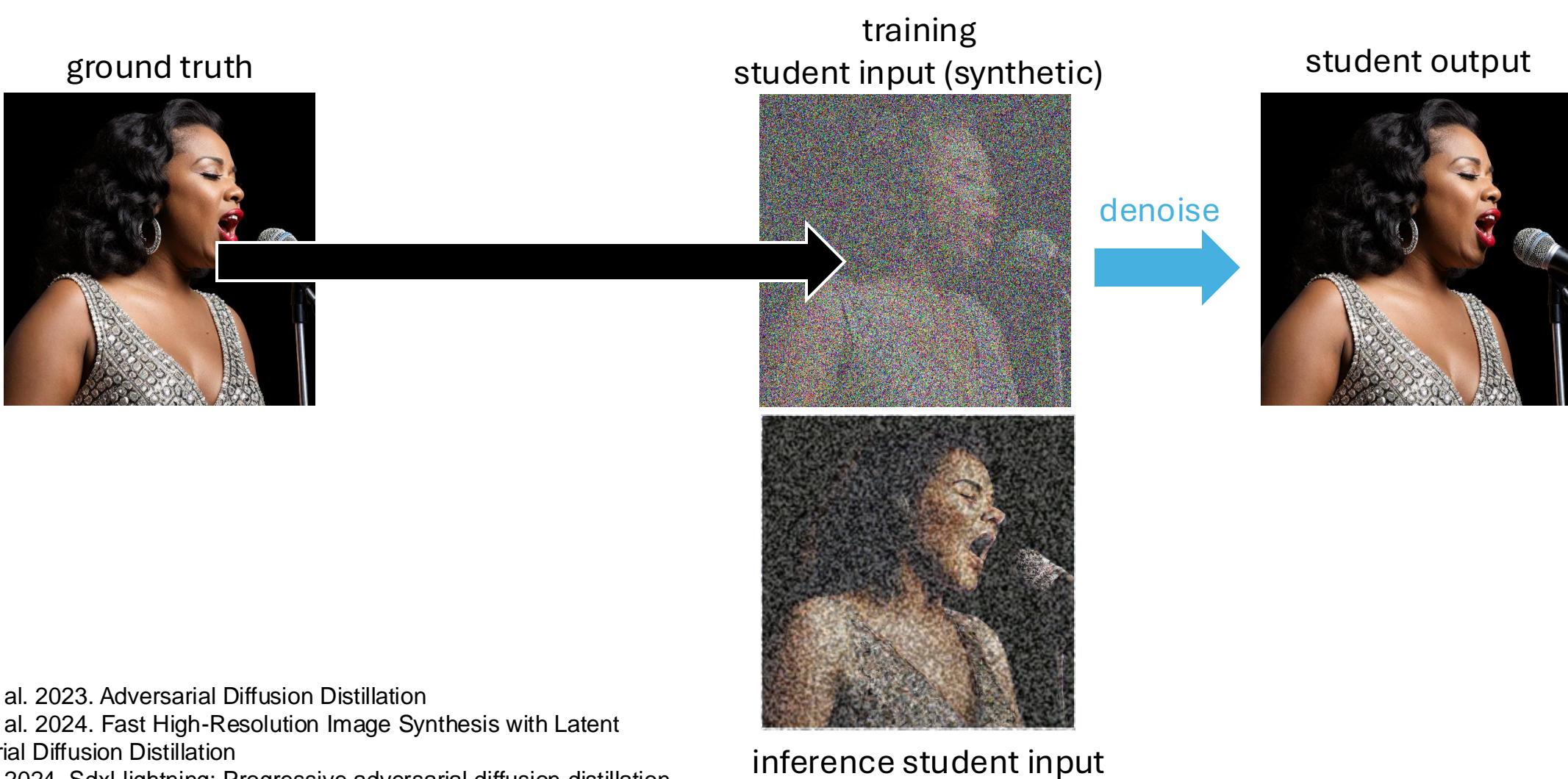
Ours 1-step (92ms)



Ours 4-step (368ms)



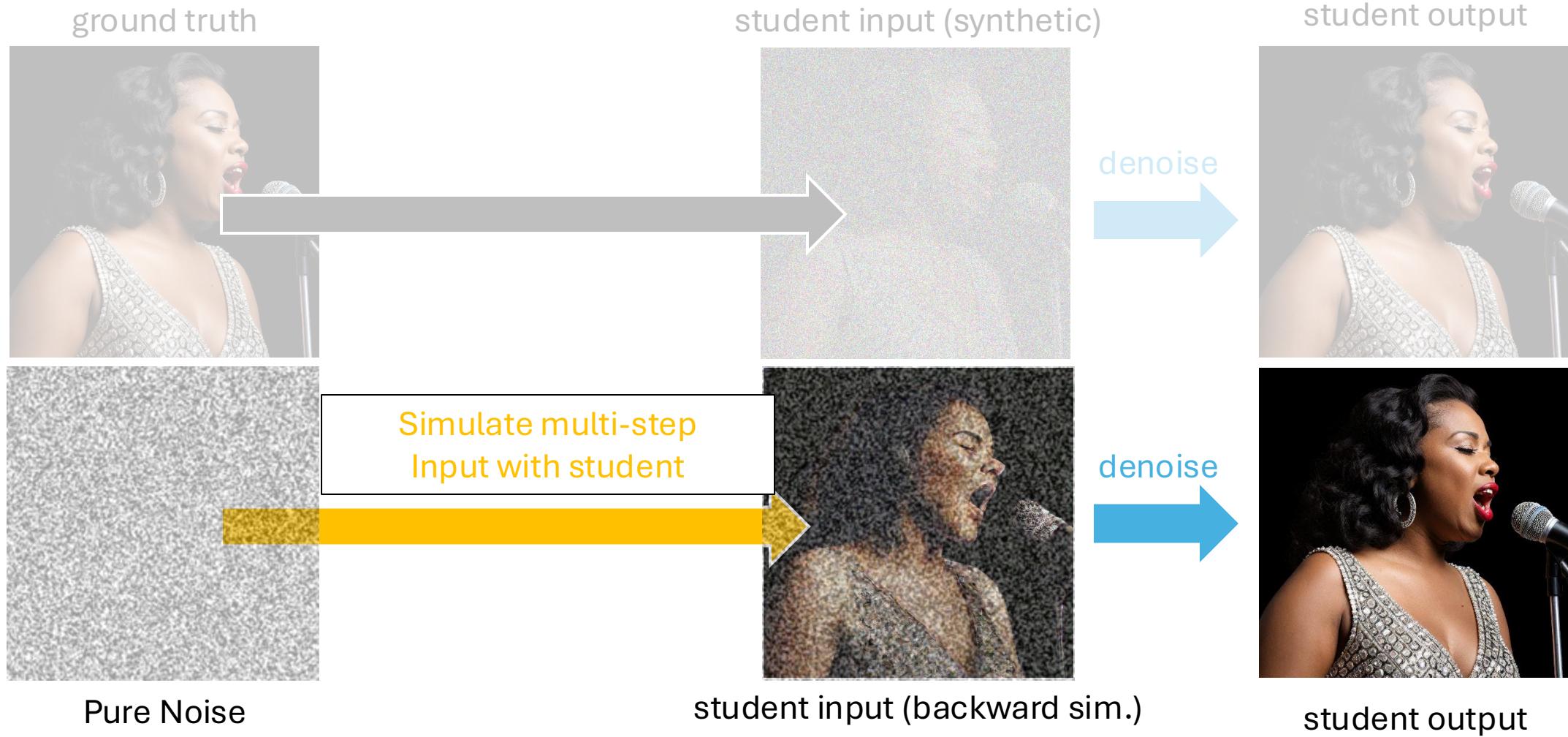
Domain gap Multi-Step Distillation



Domain gap in Multi-Step Distillation



Fix: Backward Simulation





Ours w. Backward Simulation

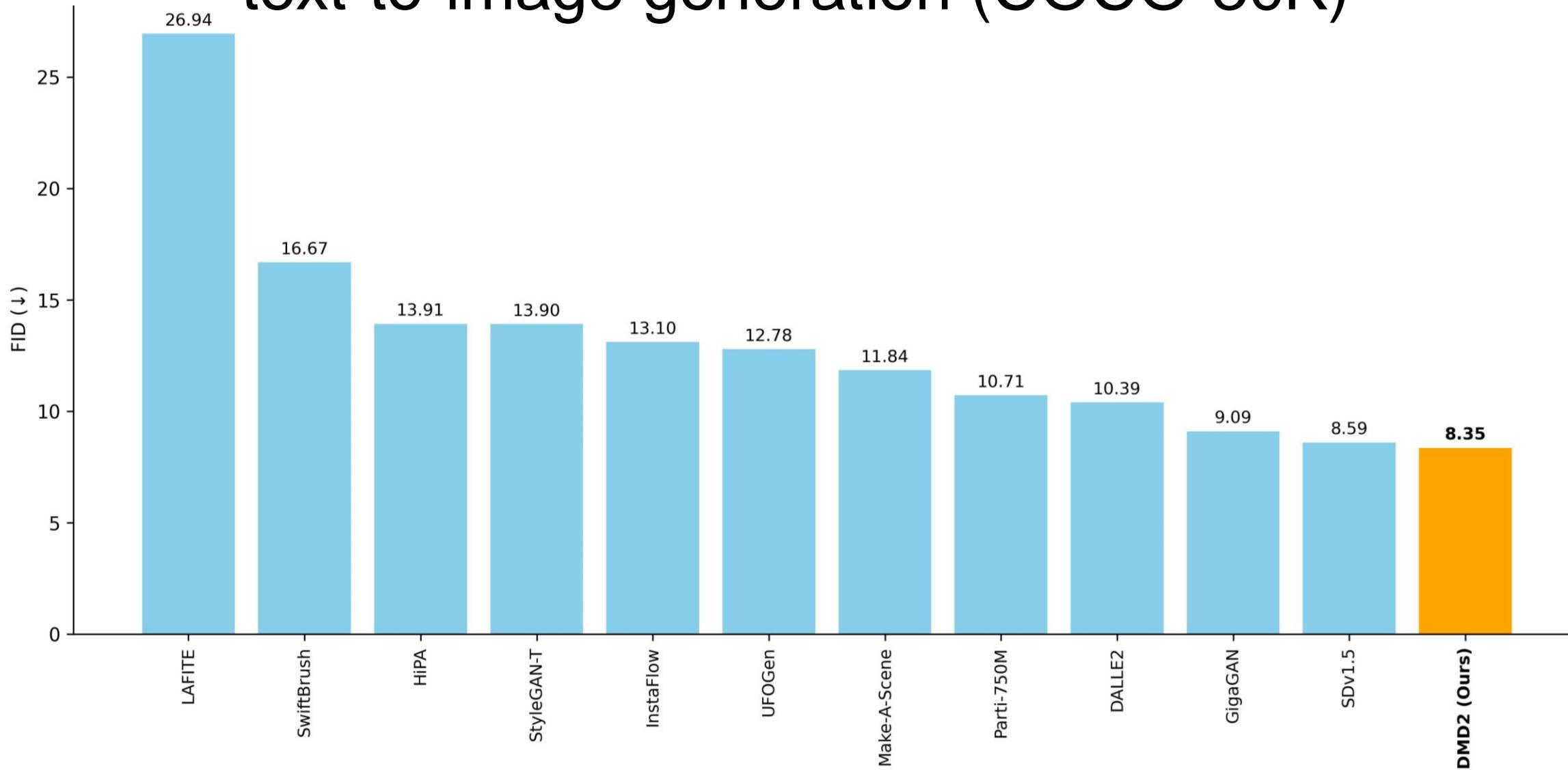


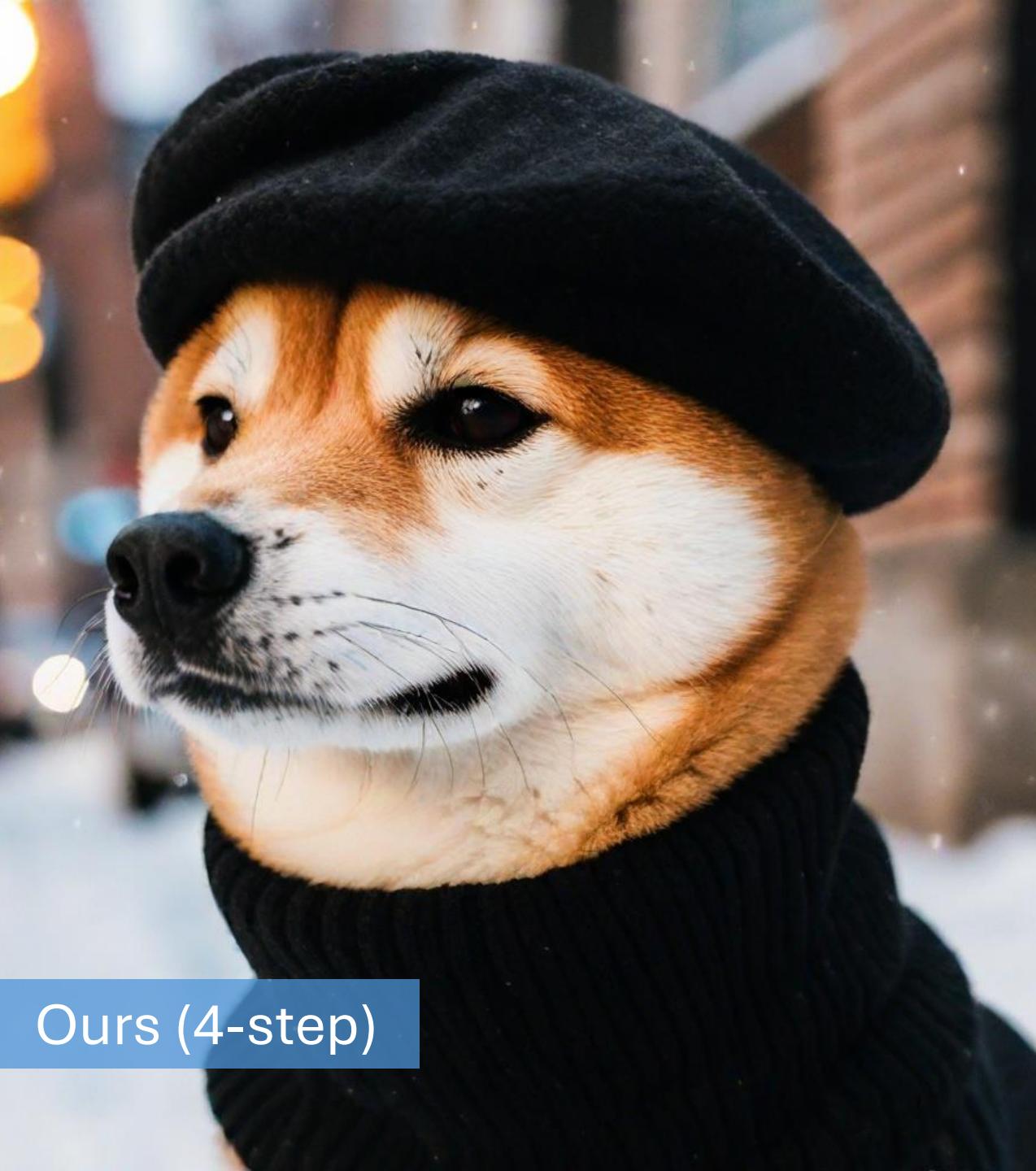
No Backward Simulation

Results

“The giant magical deer god of the forest floor, [...] analogous colors, Award-winning photography”

State-of-the-art 1-step text-to-image generation (COCO-30K)





Ours (4-step)



Teacher (SDXL, 50-step)



Ours (4-step)



SDXL-Lightning (4-step)

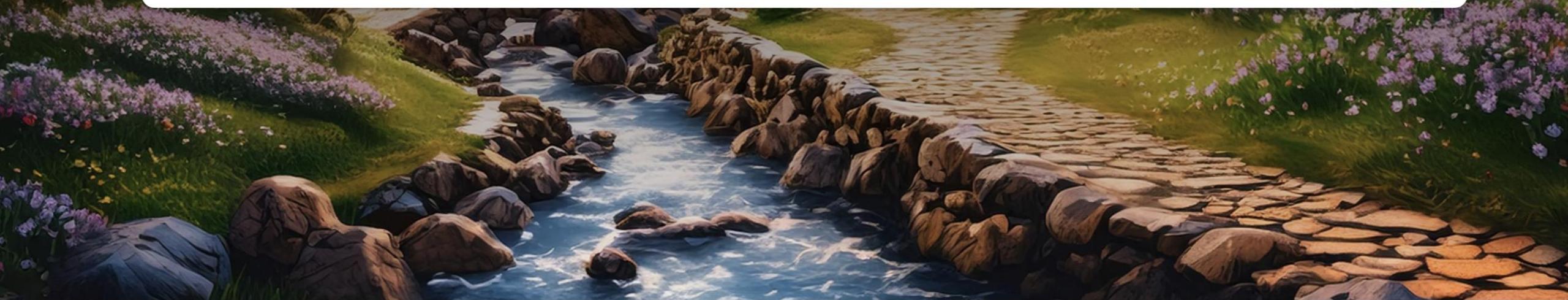


samples from a distilled 4-step generator

Get creative with Adobe Firefly

Beautiful cozy fantasy stone cottage in a spring forest aside a cobblestone path and a babbling brook. Stone wall. Mountains in the distance. Magical tone and feel, hyper realistic.

[Generate](#)



Create with generative AI

Experiment with the latest innovations from Firefly and other generative AI technology, and let us know what you think.

[Generate video](#)

Create and edit video with Firefly. Learn more and join the waitlist.

[Learn more](#)

Introducing Fast mode

Ideate, generate, and compare variations quickly. Upscale to higher resolution 2K images later for 1 credit per image.

[Try it](#)

Video Diffusion Models

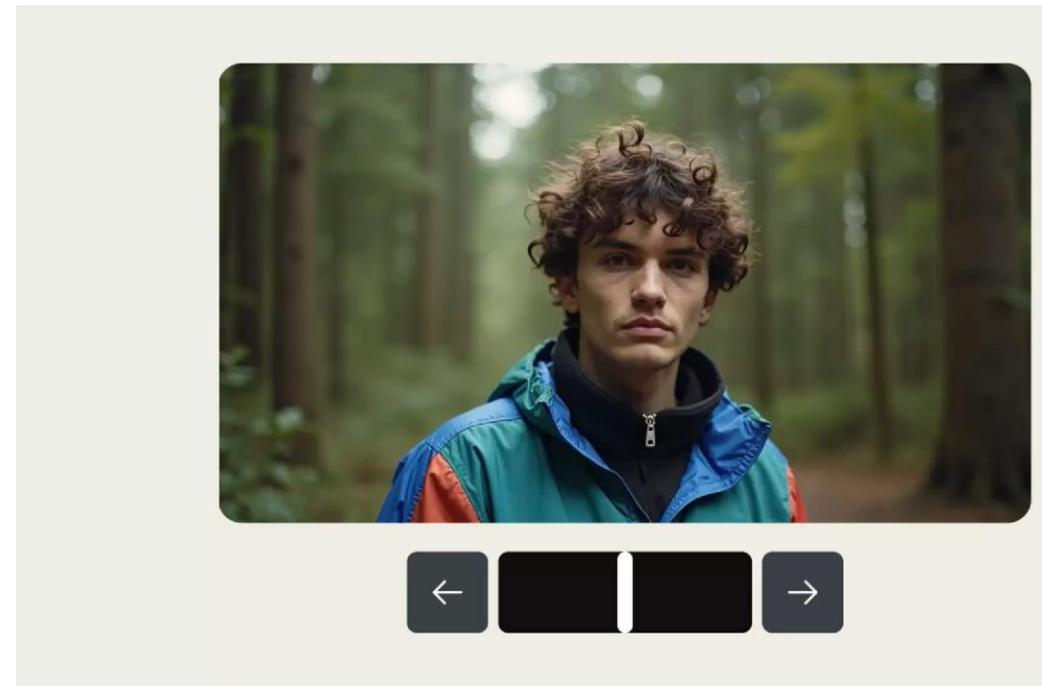


Drone view of Big Sur, generated by OpenAI's SORA



Portrait of a deer in a snowy forest, generated by Adobe's Firefly Video

Challenges in Achieving Real-Time Performance for Interactive Applications



Real-Time
Intelligence

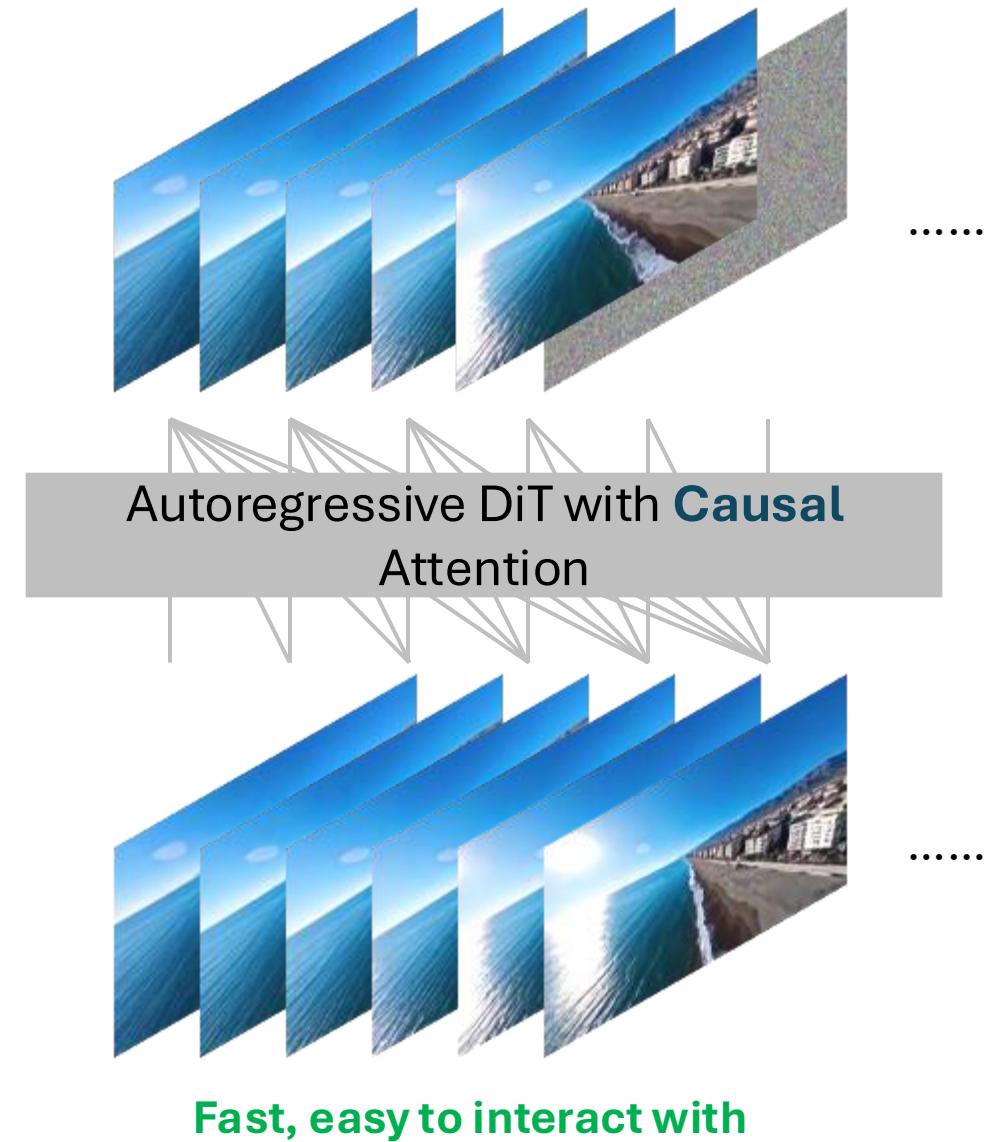
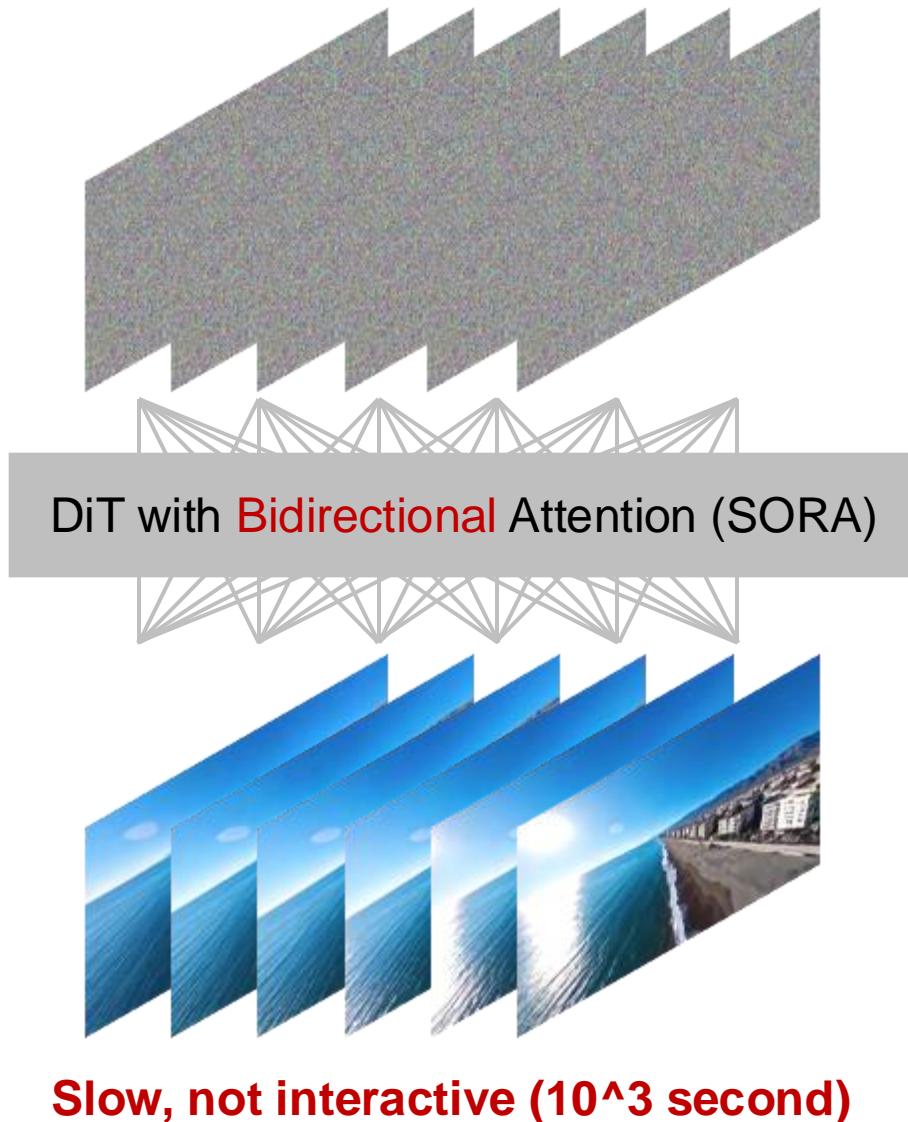
How to get the additional 10X speed up ?

on, video

Demo from Runway

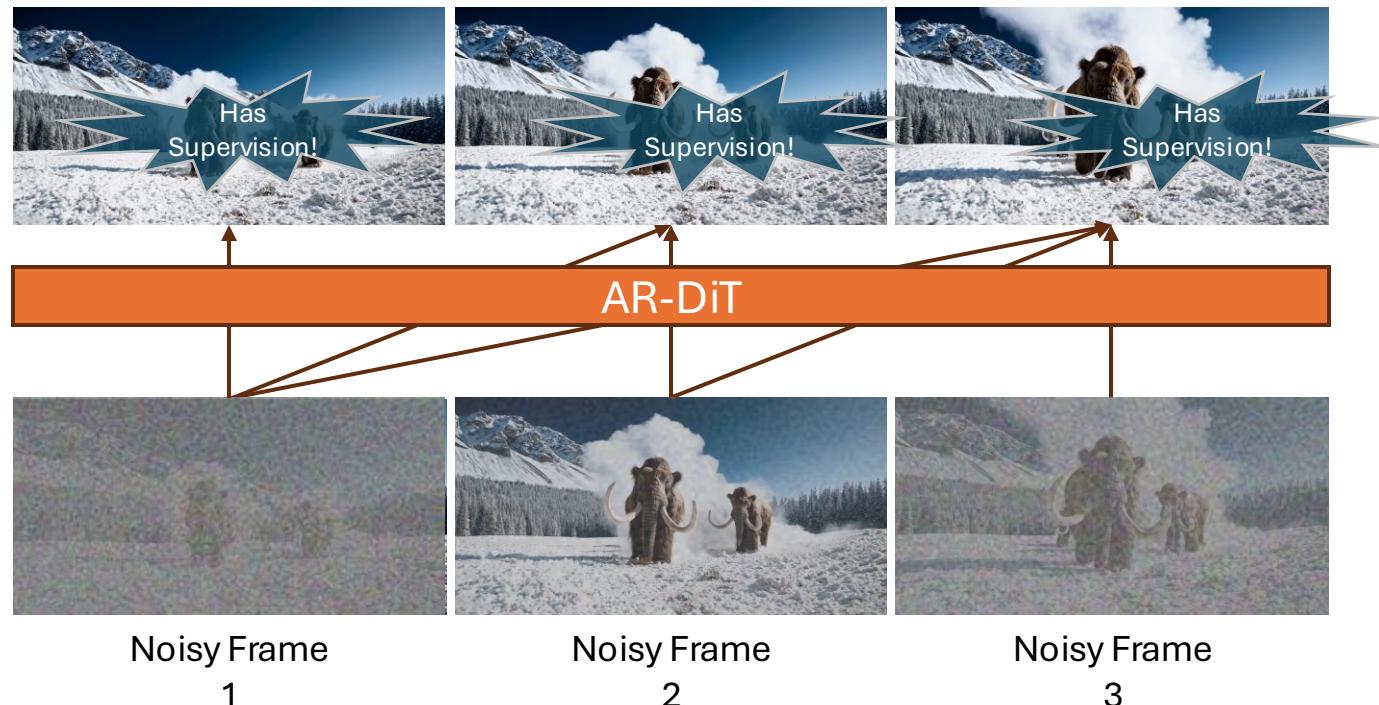
Combining Autoregressive Generation
and Distillation for Fast Video Generation

Rethinking the Bidirectional Attention in Video Diffusion Model



Scalable Decoder-only Autoregressive Diffusion Transformer

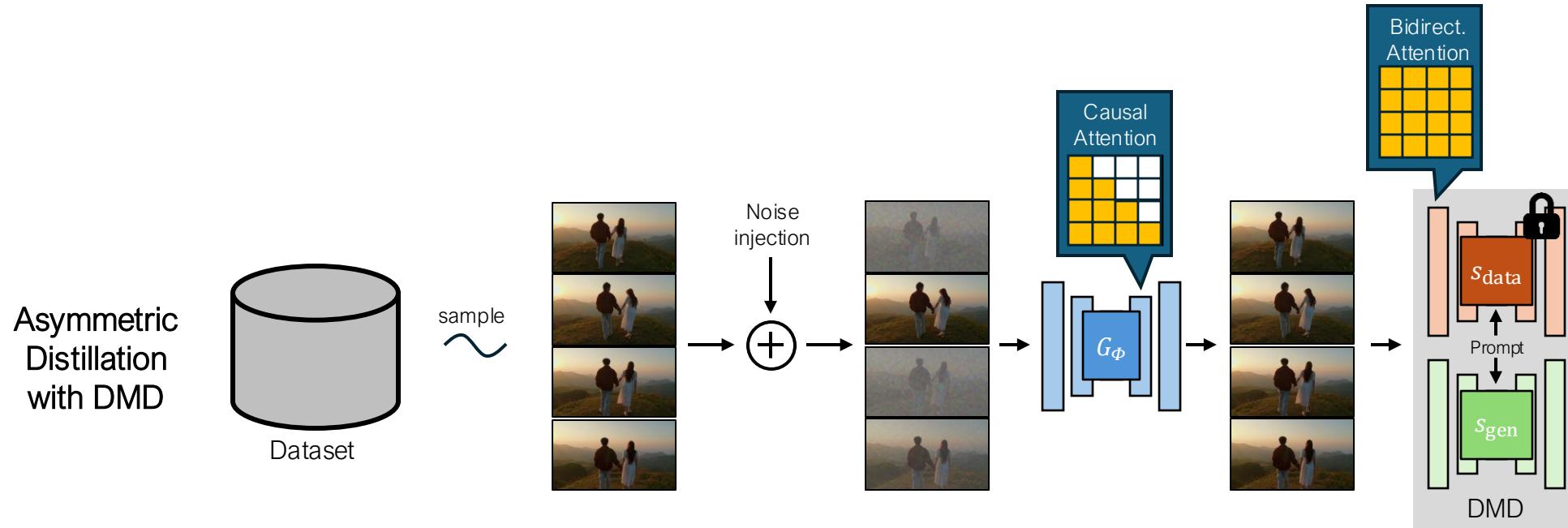
- Generalization of Diffusion-Forcing to causal DiT.
- Add noise with independent noise level to each frame.
- Given any noisy context frames, denoise all following noisy frame
- Denoising of all subsequences can be paralleled with a single forward pass using a block-wise causal mask.
- Supervision efficiency = 100%



Error Accumulation in Causal Diffusion Model



Solution: Distilling a Bidirectional Model into a Robust Causal Generator



Our Asymmetric Distillation Largely Alleviates Error Accumulation



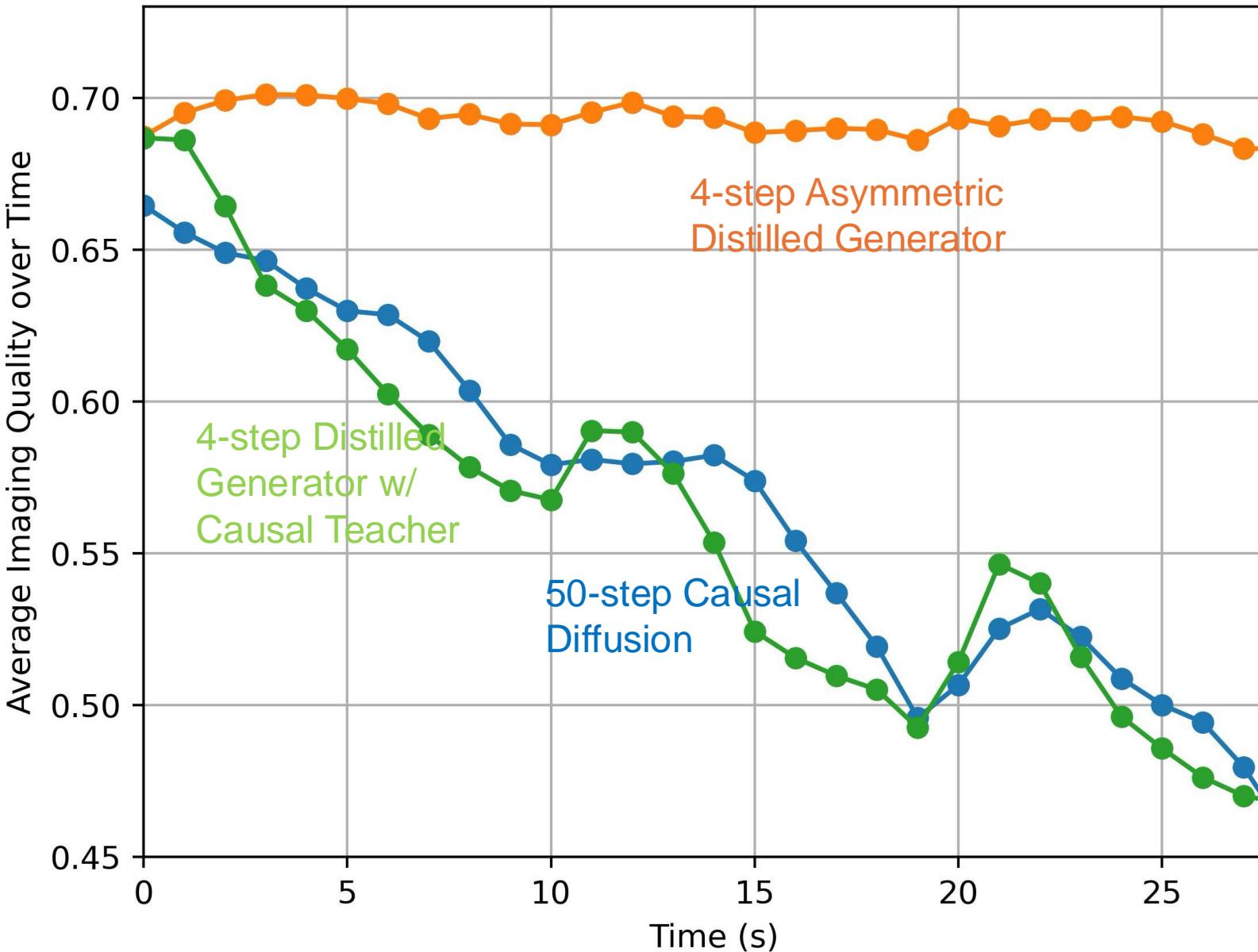
Causal Diffusion Model (50 Step)

Student Generator w/ Asymmetric Distillation (4 Step)

Distillation improves both speed and quality!

Only possible with distribution-matching style approaches!

Our Asymmetric Distillation Largely Alleviates Error Accumulation



Generators distilled from Bidirectional Teacher maintains quality over time

Intuitively, the DMD loss with a bidirectional teacher appears to be a significantly more robust objective compared to the original denoising loss.

Generators distilled from causal teacher still suffers from error accumulation

Connection with LLM: Analogy of Full Sequence-Level Distribution Matching Loss and RLHF

Pretraining

LLM

Autoregressive
Diffusion

Cross Entropy Loss

Denoising Loss

Independent for Each Token / Frame
Suffer from Error Accumulation / Hallucination

Post-Training

RLHF

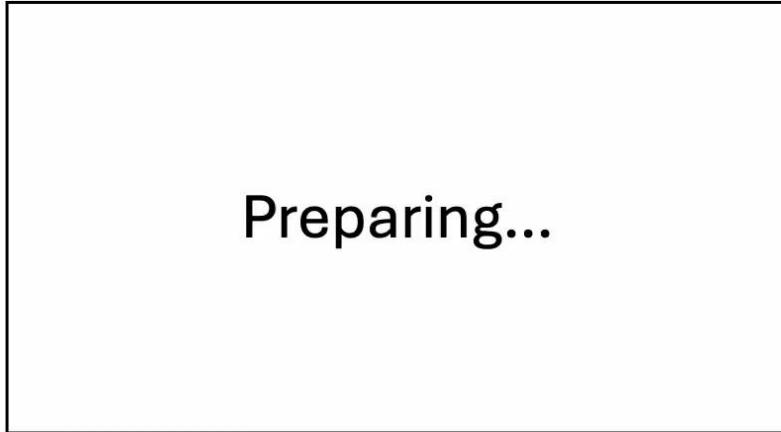
DMD with
bidirectional teacher

Full-Sequence Level Loss
Much more robust long-form rollout

Results

Enabling Real-Time Video Viewing During Generation

Bidirectional Teacher



Preparing...

Progress: 0/1



00:00

```
16] -1/-1/-1->0->-1 [17] -1/-1/-1->0->-1 [18] -1/-1/-1->0->-1 [19] -1/-1/-1->0->-1 [20] -1/-1/-1->0->-1 [21] -1/-1/-1->0->-1 [22] -1/-1/-1->0->-1 [23] -1/-1/-1->0->-1 [24] -1/-1/-1->0->-1 [25] -1/-1/-1->0->-1 [26] -1/-1/-1->0->-1 [27] -1/-1/-1->0->-1 [28] -1/-1/-1->0->-1 [29] -1/-1/-1->-1 [30] -1/-1/-1->0->-1 [31] -1/-1/-1->0->-1  
xuhuang-0464291930-0-0:504194:506058 [0] NCCL INFO P2P Chunksize set to 131072  
xuhuang-0464291930-0-0:504194:506058 [0] NCCL INFO Connected all rings  
xuhuang-0464291930-0-0:504194:506058 [0] NCCL INFO Connected all trees  
xuhuang-0464291930-0-0:504194:506058 [0] NCCL INFO 32 coll channels, 32 collnet channels, 0 nvls  
channels, 32 p2p channels, 32 p2p channels per peer  
xuhuang-0464291930-0-0:504194:506058 [0] NCCL INFO TUNER/Plugin: Failed to find ncclTunerPlugin_v2,  
using internal tuner instead.  
xuhuang-0464291930-0-0:504194:506058 [0] NCCL INFO ncclCommInitRank comm 0x55a809dbc940 rank 0  
ranks 1 cudaDev 0 nvmlDev 0 busId 53000 commId 0x54a120ee0c121fbc - Init COMPLETE
```

CausVid (Ours)



Preparing...

Progress: 0/15

```
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 5 device #3 0000:a4:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 6 device #0 0000:b8:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 6 device #1 0000:b7:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 6 device #2 0000:b6:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 6 device #3 0000:b5:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 7 device #0 0000:c9:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 7 device #1 0000:c8:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 7 device #2 0000:c7:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI NIC group 7 device #3 0000:c6:00.0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO NET/OFI Libfabric provider associates MRS with domains  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO Using non-device net plugin version 0  
xuhuang-0464291930-0-0:472752:473617 [0] NCCL INFO Using network AWS Libfabric
```

First Place on VBench Benchmark

Model Name (clickable) ▲	Sampled by	Evaluated by	Accessibility	Date	Total Score	Quality Score	Semantic Score
CausVid	CausVid Team	VBench Team		2024-12-07	83.88%	85.21%	78.57%
MiniMax-Video-01	VBench Team	VBench Team	API	2024-10-01	83.41%	84.85%	77.65%
HunyuanVideo_(Open-Source)	VBench Team	VBench Team	OpenSource	2024-12-16	83.24%	85.09%	75.82%
Gen-3_(2024-07)	VBench Team	VBench Team	API	2024-07-25	82.32%	84.11%	75.17%
Vchitect-2.0_(VEnhancer)	VBench Team	VBench Team	OpenSource	2024-09-20	82.24%	83.54%	77.06%
CogVideoX1.5-5B_(5s_SAT_prompt)	VBench Team	VBench Team	OpenSource	2024-11-15	82.17%	82.78%	79.76%
Jimeng	VBench Team	VBench Team	API	2024-11-15	81.97%	83.29%	76.69%
Vidu	VBench Team	VBench Team	API	2024-11-15	81.89%	83.85%	74.04%
Kling_(2024-07_high-performance)	VBench Team	VBench Team	API	2024-08-01	81.85%	83.39%	75.68%
CogVideoX-5B_(SAT_prompt-only)	VBench Team	VBench Team	OpenSource	2024-09-02	81.61%	82.75%	77.04%
Vchitect-2.0-2B	VBench Team	VBench Team	OpenSource	2024-09-16	81.57%	82.51%	77.79%
CogVideoX-2B_(SAT_prompt-only)	VBench Team	VBench Team	OpenSource	2024-08-19	80.91%	82.18%	75.83%
Pika-1.0_(2024-06)	VBench Team	VBench Team	API	2024-07-29	80.69%	82.92%	71.77%

Text to Video Generation

“Macro shot of a man wearing an antique diving helmet with dark glass and a jetpack walking on the veins of a leaf. Realistic style”



Text to Video Generation

“A Samoyed and a Golden Retriever dog are playfully romping through a futuristic neon city at night. The neon lights emitted from the nearby buildings glistens off of their fur.”



Text to Video Generation

“a spooky haunted mansion, with friendly jack o lanterns and ghost characters welcoming trick or treaters to the entrance, tilt shift photography.”



Long Video Generation via Sliding Window Inference

“A bear made of strawberrys is walking in the forest, its eyes looking around as if it is seeing the world for the first time.”



Zero-Shot Image to Video Generation

“Rocket blasting off from a laptop screen on an organized office table. The rocket leaves the screen and blast into space.”



Limitation: Reduced Diversity

Macro cinematography, slow motion shot: A sculptor's hands shape wet clay on a wheel, and as the wheel spins. Camera captures the tactile quality of the clay and the fluid motion of the sculptor's hands.



50-step Bidirectional Diffusion Model

4-step Causal Generator

Discussion

Trajectory Preserving

- + Better diversity
- Blurry outputs for few-step generation
- Over constrained with limited performance upper bound
- Teacher/student often requires same architecture

Mix of Both

- Hyper-SD, DMD1, SDXL-Lightning
- Explorations on standalone methods:
 - Moment Matching Distillation (DM)
 - Continuous-Time Consistency Models (TP)

Distribution Matching

- Mode collapse
- + Higher quality for few-step generation
- More unstable training
- + More flexible architecture configuration